# BRAIN-INSPIRED COMPUTATIONAL INTELLIGENCE VIA PREDICTIVE CODING

**Tommaso Salvatori**[1,2]    **Ankur Mali**[3]    **Christopher L. Buckley**[1,4]    **Thomas Lukasiewicz**[2,5],
**Rajesh P. N. Rao**[6]    **Karl Friston**[1,7]    **Alexander Ororbia**[8]

[1] VERSES AI Research Lab Los Angeles, California, USA
[2] Institute of Logic and Computation, Vienna University of Technology, Austria
[3] University of South Florida Tampa, FL 33620, USA
[4] Sussex AI Group, Department of Informatics, University of Sussex, Brighton, UK
[5] Department of Computer Science, University of Oxford, UK
[6] Paul G. Allen School of Computer Science and Engineering, University of Washington Seattle, Washington, USA
[7] The Wellcome Centre for Human Neuroimaging, Queen Square Institute of Neurology University College London
[8] The Neural Adaptive Computing Laboratory, Rochester Institute of Technology Rochester, NY 14623
`tommaso.salvatori@verses.ai, ankurarjunmali@usf.edu`
`chris.buckley@verses.ai, thomas.lukasiewicz@tuwien.ac.at`
`rao@cs.washington.edu, karl.friston@verses.ai, ago@cs.rit.edu`

## ABSTRACT

Artificial intelligence (AI) is rapidly becoming one of the key technologies of this century. The majority of results in AI thus far have been achieved using deep neural networks trained with a learning algorithm called error backpropagation, always considered biologically implausible. To this end, recent works have studied learning algorithms for deep neural networks inspired by the neurosciences. One such theory, called *predictive coding* (PC), has shown promising properties that make it potentially valuable for the machine learning community: it can model information processing in different areas of the brain, can be used in control and robotics, has a solid mathematical foundation in variational inference, and performs its computations asynchronously. Inspired by such properties, works that propose novel PC-like algorithms are starting to be present in multiple sub-fields of machine learning and AI at large. Here, we survey such efforts by first providing a broad overview of the history of PC to provide common ground for the understanding of the recent developments, then by describing current efforts and results, and concluding with a large discussion of possible implications and ways forward.

## 1 Introduction

The machine learning community is developing and producing models that push the boundaries of the field on a weekly basis: we have witnessed considerable breakthroughs in the fields of generative artificial intelligence (AI) [1, 2, 3, 4], game playing [5, 6], and text-generation [7, 8, 9]. These results reflect over a decade of advances in the field, made possible by the joint effort of tens of thousands of researchers and engineers who have built on seminal work [10, 11] in order to improve the performance of deep artificial neural networks trained with the error backpropagation algorithm [12, 13]. It may thus come as a surprise that research on alternative training methods for machine learning is more active than ever, with many works often co-authored by the same scientists who pioneered backpropagation-based schemes [14, 15, 16, 17]. These lines of research, however, are not *contradictory* but *complementary*: the increasingly impressive results of standard deep neural networks have also highlighted several significant limitations, that may be addressed via alternative training methods. There is, in fact, a common belief that, to accelerate progress in AI, we must invest in fundamental research, inspired by the findings and ideas in computational neuroscience [18, 17, 19]. To this end, in this review, we focus on a specific approach, inspired by a theory of learning and perception in the brain called *predictive coding* (PC) [20, 21, 22]. Models of PC exhibit several valuable properties, the main one being the locality of the operations, meaning that every update of the parameters is performed using only the information of adjacent neurons
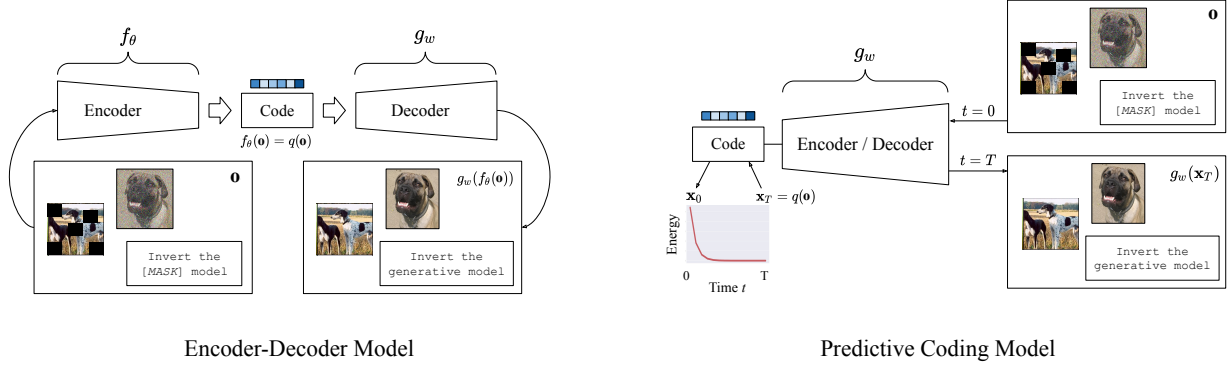
Figure 1: Generative models effectively compress information about a specific data point $\mathbf{o}$, with missing information, into a low-dimensional code vector (or latent *embedding*), and use it to generate a (e.g., semantically) similar, complete data point. Left: the standard encoder-decoder model used in machine learning [30]; Right: an equivalent PC model, which iteratively computes and refines the code through a free-energy minimization process. At convergence, the code is then used to generate a data point via the same model used to perform the compression. Typically, PC associates the encoder with ascending (prediction error) messages and the decoder with descending (prediction) messages. Note that the implicit conflation of the encoding and decoding means that there is only one set of parameters in PC. These are the parameters of the generative model that, crucially, can be optimized locally, given the requisite predictions and prediction errors.

or synapses in a Hebbian way [23]. Differently from backprop, which is a sequential algorithm, locality allows for a full parallelization of the operations, making it potentially more suitable for applications in neuromorphic computing. Locality also allows the training of artificial neural networks with any given topology [24, 25]. This facilitates the training of models with arbitrarily entangled structures, similar to the biological networks that constitute our brains and underwrite our level of intelligence [26]. Furthermore, PC networks have shown to be robust and stable due to the relaxation process that spreads error signals across the whole network, and emulates implicit gradient descent [27, 28, 29], allowing them to perform well in tasks that are more likely to be faced by natural agents.

**Notation.** In this paper, the symbol $\odot$ is employed to represent the Hadamard product, i.e., element-wise multiplication. $\oslash$ is used to indicate element-wise division. On the other hand, the symbol $\cdot$ denotes the (dot product) multiplication of matrices or vectors. The transpose of a vector $\mathbf{v}$ is represented as $(\mathbf{v})^\mathsf{T}$. It is important to note that matrices and vectors are distinguished by using boldface type, as in matrix $\mathbf{M}$ and vector $\mathbf{v}$, while scalars are depicted in italicized font, as in scalar $s$. In the context of neural networks, an input observation, or a sensory input pattern, is represented as $\mathbf{o} \in \mathbb{R}^{J_0}$, where $J_\ell$ represents the dimension of layer $\ell$ of the network, $\ell = 0$ is the input layer and $\ell = L$ is the output layer; a label is represented by $\mathbf{y} \in \mathbb{R}^{J_L}$, and a latent vector (or code) by $\mathbf{x}^\ell \in \mathbb{R}^{J_\ell}$.

## 2 Generative Models and Predictive Coding

Generative modeling provides a principled statistical framework for understanding complex data by decomposing observations into latent causes and observable effects. At its foundation, a generative model specifies a joint probability distribution over observed data $\mathbf{o}$ and latent variables $\mathbf{x}$, formally expressed as $p(\mathbf{o}, \mathbf{x}) = p(\mathbf{o} \mid \mathbf{x})p(\mathbf{x})$. This factorization separates the generative process—how latent causes produce observations through the likelihood $p(\mathbf{o} \mid \mathbf{x})$—from prior beliefs about the structure of those causes, encoded in $p(\mathbf{x})$. To build intuition, consider modeling handwritten digits: the latent variables $\mathbf{x}$ might capture high-level concepts like digit identity, writing style, or stroke thickness, while the likelihood $p(\mathbf{o} \mid \mathbf{x})$ describes how these abstract concepts combine to produce the pixel intensities we observe in an image. The prior $p(\mathbf{x})$ encodes beliefs about which combinations of these concepts are more or less likely to occur naturally.

The power of generative modeling lies in its ability to address two fundamental computational problems. The first is *inference*: given an observation $\mathbf{o}$, what are the most plausible latent causes that generated it? This requires computing the posterior distribution $p(\mathbf{x} \mid \mathbf{o})$ using Bayes' rule. The second is *learning*: given a dataset of observations, how should we adjust the model parameters to better capture the underlying data distribution? This typically involves maximizing the model evidence $\log p(\mathbf{o})$ with respect to the model parameters, integrating over all possible latent variable configurations. In practice, both inference and learning present significant computational challenges: computing the posterior $p(\mathbf{x} \mid \mathbf{o})$

requires evaluating the often intractable marginal likelihood $p(\mathbf{o})$, while learning requires optimizing the model evidence $\log p(\mathbf{o})$ over parameter space. These computational bottlenecks have motivated the development of various approximation schemes, each making different trade-offs between accuracy and computational efficiency.

**Variational Free-Energy.** The ability to memorize a training set, learn its patterns, and generalize to an unseen test set is the primary goal of any machine learning algorithm. The variational free energy provides a tractable objective that balances two critical factors: one that encourages fitting the training dataset (training *accuracy*), and one that penalizes the *complexity* of the inferred latent variables, facilitating better performance on unseen data. This enforces regularity, as it facilitates convergence towards the minimally complex model that provides an accurate fit of a specific dataset, aligned with Occam's razor formulations [31]. Given an observation $\mathbf{o}$ and an approximate posterior distribution $q_\theta(\mathbf{x})$ parameterized by $\theta$, the evidence lower bound (ELBO) or negative variational free energy is defined as:

$$\mathcal{L}(\theta, \mathbf{o}) = \underbrace{\mathbb{E}_{q_\theta(\mathbf{x})}[\log p(\mathbf{o} \mid \mathbf{x})]}_{\text{Accuracy}} - \underbrace{D_{KL}[q_\theta(\mathbf{x}) \| p(\mathbf{x})]}_{\text{Complexity}}$$

This objective provides a lower bound on the log marginal likelihood: $\log p(\mathbf{o}) \geq \mathcal{L}(\theta, \mathbf{o})$. The accuracy term encourages the model to assign high likelihood to the observed data under the inferred latent variables, while the complexity term penalizes deviations of the approximate posterior from the prior, preventing overfitting and ensuring generalization. A typical schema to maximize this bound is the expectation-maximization (EM) algorithm, where inference first finds the variational parameters $\theta$ that maximize $\mathcal{L}(\theta, \mathbf{o})$, and then the generative model parameters are updated to follow the same goal.

**Approximate Posterior.** Today's most prominent approach, the variational autoencoder (VAE) [30], uses amortized inference to compute the approximated posterior. This means that $q_\theta(\mathbf{x})$ is the output of the forward pass of a specifically trained neural network encoder, as sketched in Fig. 1 (left). However, an alternative class of methods performs inference iteratively, without relying on a separate encoder network. Predictive coding (PC) represents one such approach, using iterative optimization to approximate the posterior while maintaining the same generative model throughout both inference and learning, as sketched in Fig. 1 (right). This computational framework, while having roots in signal processing and neuroscience, provides a principled approach to generative modeling that we now describe in more detail.

## 2.1 What Defines Predictive Coding?

In the PC framework, given an observation $\mathbf{o}$, the goal is to approximate the posterior distribution $p(\mathbf{x} \mid \mathbf{o})$ through iterative optimization rather than amortized inference. Historically, this process was developed as a compression mechanism, where compression is defined as the process of assigning a lower-dimensional latent code to a particular observation. The first algorithmic formulation of PC dates back to the 1950s as a compression algorithm for time series data [32, 33, 34], and it is this *coding* that gave the framework its name. The first instances of PC in neuroscience also treated it as a pure inference algorithm, with no learning involved, as the goal was to model how spatial information was encoded in the retina [35].

A major paradigm shift occurred several years later with the work of Rao and Ballard, which employed PC as a *learning* algorithm in addition to inference, modeling hierarchical processing in the visual cortex [20]. In this formulation, synaptic weights are updated after inference to minimize the prediction error of the model. The connection with variational inference was established later, connecting PC to hierarchical Gaussian generative models [36, 37]. As a result of these independent but deeply related lines of work, a general theory has unified these views by showing that iterative updating schemes can be described as a process of minimizing a *variational free energy* [38], also known in machine learning as the evidence lower bound (ELBO) [39]. A sketch of the timeline of main breakthroughs in PC history is given in Fig. 2.

The PC literature has historically focused on a specific class of models that make two key simplifying assumptions: first, it typically assumes that all probability distributions are multivariate Gaussian; second, it often employs a hierarchical structure where latent variables follow a first-order Markov assumption—each layer depends only on the layer immediately above it. These assumptions, while restrictive, provide significant computational advantages: Gaussian distributions are conjugate to themselves under linear transformations, enabling closed-form updates, and the Markov structure allows for local computations between adjacent layers.

Let $\mathbf{o} = g(\mathbf{x}, \mathbf{w})$ be a generative model, where a vector of latent variables $\mathbf{x}$ and a set of parameters $\mathbf{w}$ are used to generate an observation $\mathbf{o}$. In our case, we consider generative models with $L$ layers, whose joint probability follows
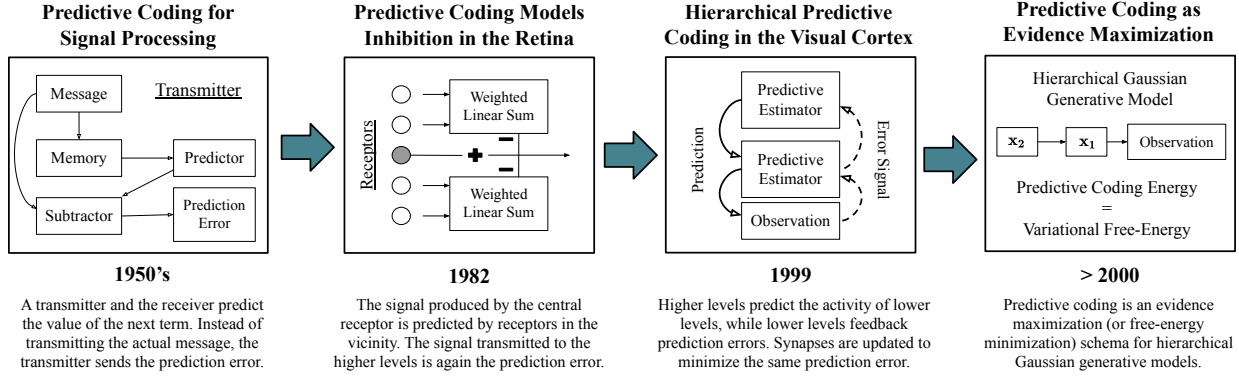
Figure 2: A timeline of how the perception of what PC is has changed through the years. Initially, it was developed as a signal compression mechanism [34, 40]; then, it was used to model inhibition in the retina [35]. It then became a more general model of both learning and perception in the visual cortex [20]. Nowadays, it can be abstractly defined as an evidence-maximization scheme for hierarchical Gaussian generative models [41, 21]. For a detailed discussion on different PC algorithms, we also refer to another survey [42]

the hierarchical dependencies:

$$p(\mathbf{x}_L, \ldots, \mathbf{x}_0) = p(\mathbf{x}_L) \prod_{\ell=0}^{L-1} p(\mathbf{x}_\ell \mid \mathbf{x}_{\ell+1}). \tag{1}$$

We further consider each conditional probability distribution $p(\mathbf{x}_\ell \mid \mathbf{x}_{\ell+1})$ to be a multivariate Gaussian distribution, whose mean is given by a transformation $g_\ell$ of the latent variables of the level above. In the majority of the literature, the map $g_\ell$ is a composition of an activation function (such as ReLU), and a transformation matrix $\mathbf{w}_\ell$ (with bias $\mathbf{b}_\ell$ possibly folded into this matrix), which results in a linear map. As a result, we then, more formally, arrive at the following:

$$p(\mathbf{x}_L) = \mathcal{N}(\mathbf{x}_L; \boldsymbol{\mu}_L, \boldsymbol{\Sigma}_L), \quad p(\mathbf{x}_\ell \mid \mathbf{x}_{\ell+1}) = \mathcal{N}(\mathbf{x}_\ell; g_{\ell+1}(\mathbf{x}_{\ell+1}), \boldsymbol{\Sigma}_\ell). \tag{2}$$

**Inversion.** Given the above, we now have the structure of the main object of study that centrally characterizes the PC literature, that is, a hierarchical generative model, which lives in a continuous state space and whose probability distributions are Gaussian in nature. However, what also defines PC is the process that is used to invert the generative model, which means estimating $p(\mathbf{x}_1, \ldots, \mathbf{x}_L \mid \mathbf{o})$ via an approximate posterior $q(\mathbf{x}_1, \ldots, \mathbf{x}_L)$. To do this, we need to appeal to two different approximation results: first, we leverage a mean-field approximation, which allows us to assume that the variational posterior factorizes into conditionally independent posteriors $q(\mathbf{x}_\ell)$; then, through the Laplace approximation, we assume that approximate posterior distributions are Gaussian in form [43]. At this point, the model can be inverted by minimizing the resulting variational free energy via gradient descent or fixed-point iterations.

As the resulting variational free energy has a quadratic form [44, 41], its gradients correspond to the linear and weighted prediction errors defined in the original computational model by Rao and Ballard [20], which was notably developed before PC was cast as variational learning and inference [41]. We have now introduced all of the core concepts needed to provide a definition of PC. This definition could be used as an umbrella for all of the variations on PC in the literature.

**Definition** (Informal). *Let us assume that we have a hierarchical generative model $g(\mathbf{x}, \mathbf{o})$, inverted using an algorithm $\mathcal{A}$. Then, $\mathcal{A}$ is a predictive coding algorithm if and only if:*

1. *it maximizes the model evidence $\log p(\mathbf{o})$ by minimizing a variational free energy,*
2. *the posterior distributions of the nodes of the hierarchical structure are factorized via a mean-field approximation, and*
3. *each posterior distribution is approximated under the Laplace approximation (i.e., random effects are Gaussian).*

Note that the above definition does not say anything explicitly about *prediction error* or properties such as *locality*, which, as mentioned earlier, are commonly used to describe PC. This is because they are not foundational to PC but rather consequences of the commitment to the aforementioned generative model: the mean field approximation enforces independence, and hence, results in locality in the update rules; the Laplace approximation simplifies variational free energy to a quadratic function, meaning that its gradients are linear prediction errors.
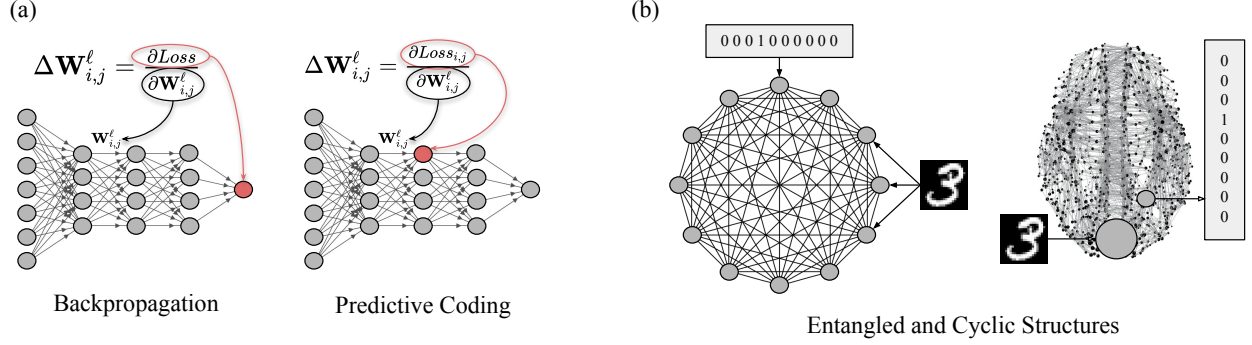
4

Figure 3: (a) The difference between PC and standard models in terms of locality: backprop updates its synaptic weights $\mathbf{W}_{i,j}^{\ell}$ to minimize the output error, even if it is not directly connected to it. PC models, on the other hand, perform their updates to correct the error of their postsynaptic neuron. (b) PC can be used to train models with cycles. An example is the fully connected model on the left, where every pair of neurons is connected via two different synapses, one in each direction. Being able to train models like this facilitates the inversion of models with arbitrarily expressive architectures (such as realistic brain structures) by simply masking specific connections of a fully connected model via an adjacency matrix. For an example of interesting models that operate on realistic brain structures, see right panel, taken from [26].

It is crucial to note that the above definition of PC does not mean, however, that the algorithms surveyed here are optimal and cannot be improved: the above definition is quite general and does not impose any constraint on the exact computation of the posteriors as well as the optimization technique(s) used to minimize the variational free energy. It also does not mean that such conditions cannot be relaxed: recent works have also explored alternatives to Gaussian and Markov assumptions. Common alternatives use amortized computations using a similar technique to VAEs [45, 46, 47, 48], rely on different probability distributions [46, 49] or use models with different architectures [25, 50]. It is also worth noting a recent trend that relies on sampling algorithms to improve the generative capacities of PC models [51, 52, 47]. However, we believe the classical framework remains essential for understanding the foundational principles that give PC its characteristic properties.

## 3 Implementations of Predictive Coding

Informally, PC claims that there exist two families of neurons that define an internal (i.e., generative) model of the world: the first generates predictions that are passed to lower layers, and the second encodes prediction errors that are passed to higher layers [20, 53]. Recent work has applied this basic separation to prediction and prediction errors within the compartments of a single neuron, where errors are propagated back via dendritic connections [54, 55]. In this work, we use the standard formulation of PC papers in machine learning, and describe an artificial neuron as a computational unit with three quantities: its value (*value node*), its *prediction*, and its error (*error node*), defined as the difference of the first two. The *value node* $\mathbf{x}_{i,t}^{\ell}$ (where the indices $i$ and $\ell$ refer to the $i$-th neuron of the $\ell$-th layer at time $t$) encodes the most likely value of some latent state. The second computational unit is the *prediction* $\mathbf{u}_{i,t}^{\ell}$, which is a function of the value nodes of a higher level in the hierarchy. Let $J_\ell$ denote the dimension (number of neurons) of layer $\ell$, $\mathbf{W}^{\ell+1}$ be a matrix containing the predictive synaptic parameters for layer $\ell + 1$, and $\phi$ be a nonlinear activation function. Then, the prediction is computed as follows:

$$\mathbf{u}_{i,t}^{\ell} = \sum_{j=1}^{J_\ell} \mathbf{W}_{i,j}^{\ell+1} \phi(\mathbf{x}_{j,t}^{\ell+1}), \tag{3}$$

where $\phi$ is a nonlinear function (akin to the activation function in deep neural networks). The third computational unit is the prediction error $\mathbf{e}_{i,t}^{\ell}$, given by the difference between its value node and its prediction node, i.e., $\mathbf{e}_{i,t}^{\ell} = \mathbf{x}_{i,t}^{\ell} - \mathbf{u}_{i,t}^{\ell}$. This local definition of error, which exists in every network neuron, foregrounds a key difference between PC and models trained with backprop (e.g., the multilayer perceptron), as it enables learning through only local computations. A graphical example depicting the differences in locality between PC and backprop is given in Fig. 3. Taken together, the above three quantities, as well as the set of synaptic weight matrices $(\mathbf{W}^0, \dots, \mathbf{W}^L)$, define a generative model where both inference and learning are performed as a means to minimize a single (global) energy function, formally

**Algorithm 1** Updating the value/state dynamics and synaptic weights in a PC model, given a data point $\mathbf{o}$.

**Require:** $\mathbf{x}^L$ is fixed to $\mathbf{o}$ for every time step $t$.
 1: **for** $t = 0$ to $T$ (included) **do**
 2:    **for** each neuron $i$ in each level $\ell$ **do**
 3:       Update $\mathbf{x}_{i,t}^\ell$ to minimize $\mathcal{E}_t$ via Eq. (5)
 4:    **if** $t = T$ **then**
 5:       Update $\mathbf{W}^\ell$ to minimize $\mathcal{E}_t$ via Eq. (6)

defined as the sum of the squared prediction errors of every neuron:

$$\mathcal{E}_t = \frac{1}{2} \sum_{i,\ell} (\mathbf{e}_{i,t}^\ell)^2. \tag{4}$$

This energy function is exactly the variational free energy defined in the previous section, and a proper derivation can be found in [56]. Let us assume that our generative model is presented with an observation $\mathbf{o} \in \mathbb{R}^{J_0}$. Then, the following process describes the credit assignment and consequent update of the synaptic parameters of the model: first, the neurons of the lowest layer are set equal to the sensory observation, i.e., $\mathbf{x}^0 = \mathbf{o}$. Next, the unconstrained neural activities are updated until the convergence step $T$ (or, for a fixed number $T$ of iterations) to minimize the energy of Eq. 4 via gradient descent. In particular, the equation for the update dynamics of the value nodes is the following:

$$\Delta \mathbf{x}_t^\ell = -\gamma \cdot \frac{\partial \mathcal{E}_t}{\partial \mathbf{x}_t^\ell} = \begin{cases} \gamma \cdot (\phi'(\mathbf{x}_t^\ell) \odot (\mathbf{W}^\ell)^\mathsf{T} \cdot \mathbf{e}_t^{\ell-1}) & \text{if } \ell = L \text{ (and unclamped)} \\ \gamma \cdot (-\mathbf{e}_t^\ell + \phi'(\mathbf{x}_t^\ell) \odot (\mathbf{W}^\ell)^\mathsf{T} \cdot \mathbf{e}_t^{\ell-1}) & \text{if } 0 < \ell < L \\ \gamma \cdot (-\mathbf{e}_t^\ell) & \text{if } \ell = 0 \text{ (and unclamped)}, \end{cases} \tag{5}$$

where $\gamma$ is the learning rate of the neural activities (usually read as a precision or inverse variance). This optimization process drives the underlying credit assignment process of PC and computes the best configuration of value nodes needed to perform a synaptic weight update. When this process has converged, the value nodes are then frozen and a single weight update is performed (via gradient descent) to further minimize the same energy function:

$$\Delta \mathbf{W}^{\ell+1} = -\alpha \left( \partial \mathcal{E}_T / \partial \mathbf{W}^{\ell+1} \right) = \alpha \left( \mathbf{e}_T^\ell \cdot \phi(\mathbf{x}_T^{\ell+1})^\mathsf{T} \right), \tag{6}$$

where $\alpha$ is the learning rate for the synaptic update. The alternation of these two phases, i.e., the value node update and weight update steps, defines the learning algorithm used to train PC networks [20, 53]. Importantly, although every computation is local, both update rules minimize the same energy function, which is globally defined over the entire network, as described in 1. The difference in locality between a weight update of backprop and the one for PC is given in Fig. 3(a).

**Supervised Learning.** In the above, we have described how to use PC to perform unsupervised learning tasks. However, this algorithm or variations of it have been shown to obtain a performance comparable to deep networks learned with backprop on supervised learning tasks [53, 57, 58]. To extend the above formulation to the domain of labeled data, we have to re-imagine the PC network as *generating the label* $\mathbf{y}$ and that the data point $\mathbf{o}$ serves as a prior on the value of the neurons in the first layer of the network. Let us assume that we have a labeled data point $(\mathbf{o}, \mathbf{y})$, where, again, $\mathbf{y}$ is the label [1]. Supervised learning is then realized by fixing the value nodes of the first and last layer to the entries of the data point and its label, respectively, i.e., $\mathbf{x}_{0,t} = \mathbf{o}$ and $\mathbf{x}_{L,t} = \mathbf{y}$ for every time step $t$. A graphical depiction of the difference between supervised and unsupervised learning using PC is given in Fig. 4.

**Generalization to Arbitrary Topologies.** It is often remarked that the networks in the brain are not strictly hierarchical [24, 60] but are extremely entangled and full of cyclic connections. Recent work has shown how PC can be utilized to model networks with any kind of structure, making it ideal to digitally perform learning tasks that require brain-like architectures such as parallel cortical columns or sparsely connected brain regions [25]. To do this, the more general concept of a PC graph is required. Let $G = (V, E)$ be a fully connected directed graph, where $V$ is a set of $n$ vertices $\{1, 2, \ldots, n\}$ and $E \subseteq V \times V$ is a set of directed edges between them, where every edge $(i, j) \in E$ has a weight parameter encoded in an $n \times n$ weight matrix $\mathbf{W}$. As a result, the predictions come from all of the neurons of the network, and the free energy of the model is again the summation over all of the (squared) prediction errors:

$$\mathbf{u}_{i,t} = \sum_{j=0}^{n} \mathbf{W}_{j,i} \phi(\mathbf{x}_{j,t}) \quad \text{and} \quad \mathcal{E}_t = \frac{1}{2} \sum_{i=0}^{n} (\mathbf{e}_{i,t})^2. \tag{7}$$

---

[1] In supervised learning, the total energy of the model can also be decomposed as $\mathcal{E} = \tilde{\mathcal{E}}_t + \mathcal{L}$, with $\tilde{\mathcal{E}}_t$ be the total energy of all the internal neurons only, and $\mathcal{L}$ be the squared error defined on the label, as in standard regression tasks. This allows to define interesting similarities with backprop [59].
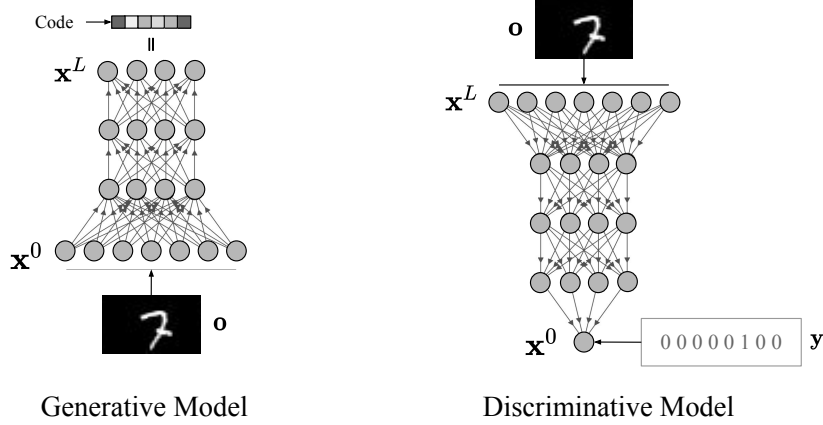
Figure 4: A graphical depiction of PC in an unsupervised generative form (left) and in a supervised discriminative form (right). Note that the generative form of PC entails iteratively inferring a latent code (or embedding) $\mathbf{x}^L$ for sensory input $\mathbf{x}^0 = \mathbf{o}$, while the discriminative form of PC requires iteratively learning a predictive mapping between sensory input $\mathbf{x}^L = \mathbf{o}$ to target $\mathbf{x}^0 = \mathbf{y}$.

While defining this generalization on fully connected graphs makes the notation simpler, interesting use cases can be found in graphs that are not fully connected. To this end, it suffices to note that every graph is a subset of a fully connected one and, hence, we can consider sparse weight matrices $\mathbf{W}$, where only the entries of the edges that we need to define our structure are non-zero. This can be done by multiplying $\mathbf{W}$ by an adjacency matrix. In the original work along these lines, the authors trained a network of multiple, sparsely connected brain regions and tested the resulting model on image denoising and image completion tasks. The exact structure of this network is that of the *Assembly Calculus*, a Hebbian learning framework/method specifically designed to model brain regions [61]. A graphical sketch of how to derive such a model is given in Fig. 3(b).

**Associative Memories.** A body of literature has investigated the efficacy of PC networks in memory-related tasks. Specifically, this literature studies their capacity to store a dataset and retrieve individual data points when a related cue is presented. In practice, this cue is often an incomplete or corrupted version of a stored memory, and the task is successfully completed if and only if the correct data point is retrieved. PC schemes have demonstrated their ability to effectively store and retrieve complex memories, such as images from the COIL and ImageNet datasets [62, 63, 64]. However, while their retrieval capabilities are robust, their capacity is still limited, and hence not comparable to that of continuous or universal Hopfield networks [65, 66]. This capacity can be improved by implementing a fast and powerful memory writing operation that allows these models to store individual memories without overwriting existing ones. Models that incorporate memory writing also tend to facilitate the implementation of a *forget* operation that erases individual memories with little to no impact on overall model performance [67]. It is also possible to store memories in the latent variable $\mathbf{x}_L$, by learning a Gaussian mixture prior, whose centers correspond to each of the individual memories [68]. This has also been implemented in the temporal domain, where the goal is to retrieve missing future frames of a video, given initial ones [64] (as a sort of priming stimulus). Note, however, that this class of models does not have the exponential capacity of modern Hopfield networks [66, 69].

### 3.1 Differences and Similarities with Backpropagation

In the case of supervised learning, PC can be seen simply as a training algorithm for deep neural networks that can be used as an alternative to backprop. However, the two algorithms differ in terms of convergence, complexity, and performance, in ways that we now discuss.

**Similarities.** Recent work has shown that PC can approximate the weight update of backprop, under specific conditions, on both MLPs (multilayer perceptrons) and within the more general framework of computational graphs [53, 44, 70, 71]. These conditions are restrictive in practice, as these results only hold if either the total prediction error on the network is infinitesimally small or the predictions are kept constant throughout the whole duration of the inference process, i.e., $\mathbf{u}_{i,t}^\ell = \mathbf{u}_{i,0}^\ell$ for every time step $t$. In practice, however, empirical studies have shown that for the approximation to hold, it suffices to have a small output error [44]. A similar line of work has also shown that simply adding a temporal scheduling to the updates of the weights and the neural activities leads to a weight update

that is equivalent to that of backprop [72, 73]. This temporal scheduling, that sequentially updates the weights of different layers at different, pre-defined, time steps, is however implausible, as it requires external control that triggers the updates at different time steps [74].

**Differences.** PC has been shown to perform better than standard models (such as backprop-trained deep neural networks) on problems that are faced by biological organisms, such as continual learning, online learning, and learning from a small amount of data [75, 76, 27]. This is due to the inference phase, which allows the error to be distributed in the network in a way that avoids a phenomenon called *weight interference*. A second difference is about stability and convergence: in backprop, and in models trained with gradient descent on a non-local loss function, each parameter update is computed independently based on the current network state, without accounting for how other parameters will simultaneously change. This has been shown to produce updates that destabilize training [77]. This is not true for PC networks, as it has been shown that PC models trained for supervised learning naturally implement *implicit* gradient descent [28], a more stable optimization approach where each parameter update takes into account the simultaneous changes that will occur to all other parameters in the network. Specifically, a parameter update of implicit gradient descent is defined as follows:

$$\mathbf{W}_{t+1} = \mathbf{W}_t - \gamma \nabla \mathcal{L}(\mathbf{W}_{t+1}). \tag{8}$$

This formulation is called *implicit*, since $\mathbf{W}_{t+1}$ appears on both sides of the equation, desirably reducing the sensitivity to the learning rate [78]. As a consequence, PC models tend to be more robust and better calibrated than standard ones, as has been shown in multiple classification tasks on both convolutional and graph networks [60, 79, 57]. Recent works have also studied and analyzed the convergence of PC models, showing that the iterative inference phase allows them to better escape bad local minima and saddle points, as well as needing a smaller number of weight updates with respect to backpropagation [29, 80].

**Limitations.** Differently from backprop-based models, PCNs face a scalability bottleneck that limits their effectiveness on large-scale deep learning tasks, preventing any kind of industrial adoption. For example, in classification tasks, these models achieve competitive performance when tested on standard benchmarks with shallow architectures, such as convolutional models with 5-7 layers. However, we observe that performance degrades significantly as network depth increases [81]. This has been shown both in PC models trained with gradient descent, and in models trained via closed-form updates [82]. In both cases, the core issue stems from energy propagation dynamics: deeper networks experience exponentially larger energy magnitudes in layers closer to the output, creating unstable training conditions where parameter updates concentrate in a small subset of layers rather than utilizing the full network depth. This energy imbalance causes latent states to diverge excessively from forward pass activations, leading to suboptimal weight updates and poor credit assignment throughout the network. While strategies such as energy regularization, modified weight update schemes, and initialization strategies are improving the results [83, 84, 85], there is still a significant amount of research to be performed to be able to make PC match the performance of BP on larger architectures, such as language models [86].

A second limitation concerns the computational complexity of the inference phase. Theoretically, optimal performance requires waiting until convergence before updating weights, which introduces substantial computational overhead. In practice, however, a limited number of inference iterations yields competitive results, typically requiring iterations at least equal to the network depth, with empirical evidence suggesting that optimal performance occurs when $L < T < 2L$ iterations are performed [81]. When the iteration count approximates the layer count, the time complexity becomes comparable to backprop, though space complexity remains higher due to the need to maintain all neurons and layers simultaneously in memory. This computational profile creates significant engineering challenges, particularly the absence of deep learning libraries capable of fully parallelizing the inference phase. Current implementations fall into two distinct categories: libraries that achieve full parallelization but lack seamless integration with established deep learning frameworks like PyTorch and Equinox, such as NGC-Learn [24] and jPC [87]; and libraries that provide modular, framework-compatible approaches but sacrifice parallelization, such as PCX [81].

## 3.2 Neural Generative Coding

Despite being developed to model information processing in different brain areas, the above formulation is still biologically implausible, the main implausibility being symmetric weight connections. In what follows, we survey a generalization of PC, to arbitrary wiring patterns, called neural generative coding (NGC). This variation has its roots grounded in cable theory [88, 89] and neuronal compartments [90, 88, 89]. Like PC, NGC employs predict-then-correct learning but incorporates additional neurobiological mechanisms including lateral competition, learned precision-weighting, and membrane potential leakage. The key innovation of NGC is decoupling the forward generative pathway (matrices $\mathbf{W}^\ell$) from the feedback pathway (matrices $\mathbf{E}^\ell$). This allows flexible connectivity patterns, including
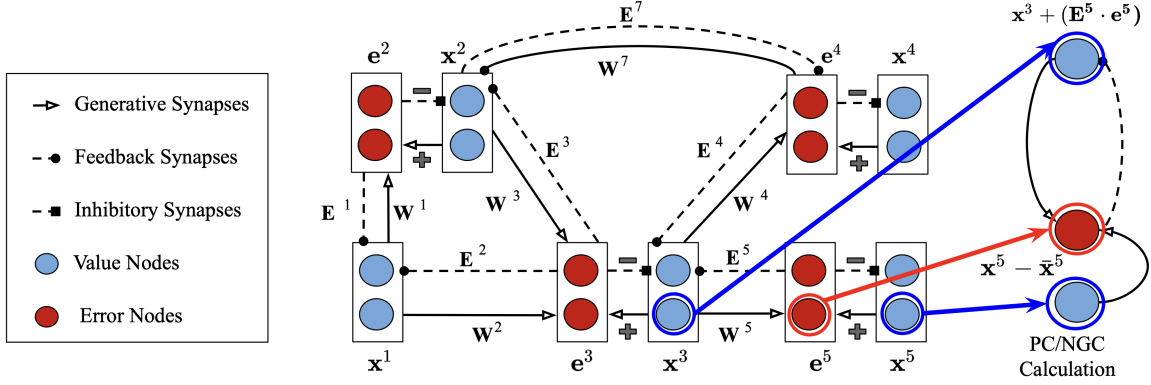
Figure 5: On the left is an arbitrary neural circuit (non-linearities are omitted for simplicity) that could be constructed under one of the three PC frameworks surveyed. Rao & Ballard PC (PC) only employs generative synapses $\mathbf{W}^\ell$, while NGC further employs separate feedback synapses (which may or may not correspond symmetrically to the generative pathway). Note that NGC/BC-DIM notably often employ skip connections ($\mathbf{W}^7$).

skip connections, without requiring that feedback structures mirror the generative model in reverse [24]. Each NGC circuit processes input over $T$ time steps, where $g^\ell$ represents a precision-weighting function, $\beta$ controls the update rate, $\gamma$ represents membrane leakage, $f_D$ is a dendritic processing function, $\Delta t$ is the time step size, and $\Phi(\mathbf{x}_t^\ell)$ captures lateral interactions. Predictions are computed and neural activities are updated according to:

$$\bar{\mathbf{x}}_t^\ell = g^\ell\Big(\mathbf{W}^{\ell+1} \cdot \phi^{\ell+1}(\mathbf{x}_t^{\ell+1})\Big), \tag{9}$$

$$\mathbf{x}_{t+\Delta t}^\ell = \mathbf{x}_t^\ell + \beta\Big(-\gamma\mathbf{x}_t^\ell - \mathbf{e}_t^\ell + (\mathbf{E}^\ell \cdot \mathbf{e}_t^{\ell-1}) \odot f_D(\mathbf{x}_t^\ell) + \Phi(\mathbf{x}_t^\ell)\Big). \tag{10}$$

Error signals are calculated as precision-weighted differences between actual and predicted activities, modulated by cross-correlational matrices $\mathbf{\Sigma}^\ell$. After settling for $T$ steps, both generative and feedback synapses are updated via multi-factor Hebbian rules with $\gamma_e < 1$ ensuring feedback synapses evolve more slowly than generative ones:

$$\Delta\mathbf{W}^\ell = \mathbf{e}_t^\ell \cdot (\phi^{\ell+1}(\mathbf{x}_t^{\ell+1}))^T, \tag{11}$$

$$\Delta\mathbf{E}^\ell = \gamma_e(\phi^{\ell+1}(\mathbf{x}_t^{\ell+1}) \cdot (\mathbf{e}_t^\ell)^T). \tag{12}$$

For a detailed explanation of all the mechanisms of NGC, we refer to the supplementary material.

## 4 Predictive Coding in Machine Learning

Despite the limitations regarding the scale highlighted in the previous chapter, PC has demonstrated remarkable versatility across a large number of machine learning tasks. While early implementations focused primarily on simple classification tasks, recent advances have shown that PC can effectively tackle increasingly sophisticated problems that were previously dominated by backpropagation-based approaches. These developments that we now survey, highlight both the practical potential and current limitations of predictive coding as a comprehensive learning paradigm for modern machine learning applications.

**Supervised Learning.** The first application of PC to supervised learning involved training a small PC network to perform image classification on the MNIST dataset, achieving test and train errors comparable to those of an MLP of the same complexity (depth and width) [53]. Since then, similar results have been achieved on convolutional networks trained on datasets of *RGB* images, such as CIFAR10 and SVHN [57]. In this work, the authors showed that updating the weight parameters alongside the neural activities (i.e., running Equations 5 and 6 in parallel), instead of waiting for the inference phase to converge, improves test accuracies as well as ensures better convergence. These results also extend to structured datasets and graph neural networks, where PC is again able to match the performance of backprop on different benchmarks, with the advantage of learning models that are better calibrated and more robust to adversarial examples [60].

**Natural Language Processing.** Gaussian assumptions (that underwrite PC) can be limiting in scenarios where we need to model different distributions, such as categorical distributions or mixture models. One such scenario is found

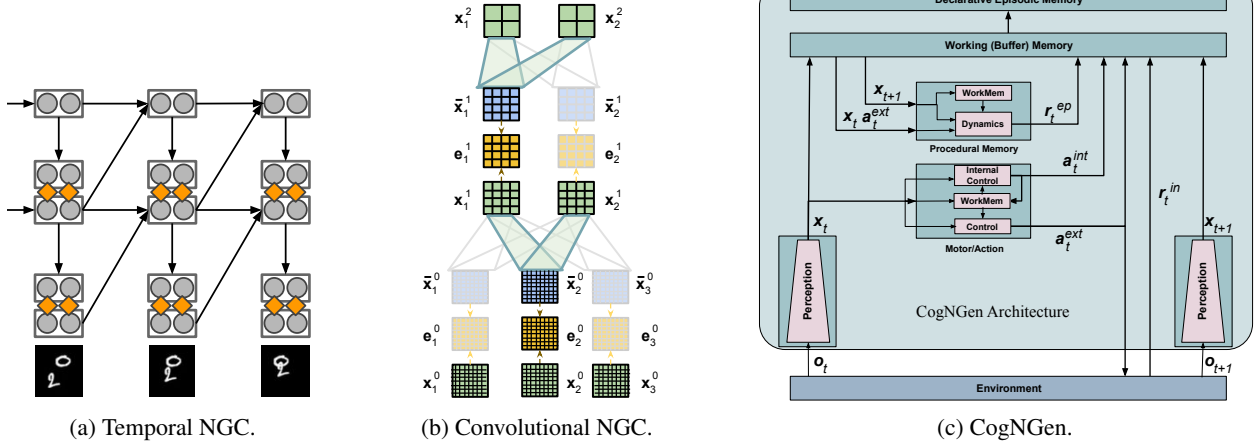(a) Temporal NGC.          (b) Convolutional NGC.          (c) CogNGen.

Figure 6: Various systems based on NGC-based PC: (a) a temporal NGC model (orange diamonds depict error nodes, grey circles depict state nodes, while arrows represent synaptic projections) for processing temporally-varying sequences of data patterns, e.g., frames from the moving MNIST dataset, (b) a convolutional NGC model (figure adapted from [93]; green grid blocks represent state feature maps, blue grid blocks represent prediction maps, yellow-orange grid blocks represent error maps, and blue parallelograms depict deconvolutional pathways) for extracting distributed representations of natural images, and (c) the CogNGen architecture (figure adapted from [94]), a modular control system that implements several core elements of the Common Model of Cognition [95] with PC and vector-symbolic memory ($\mathbf{r}$ depicts a reward, while $\mathbf{a}$ depicts an action vector).

in transformer models [91]: the attention mechanism encodes a categorical distribution, computed via the softmax activation. To this end, it is possible to generalize the definition of the energy of a specific layer, originally defined as the squared distance between the predictions $\mathbf{u}^\ell$ and the activities $\mathbf{x}^\ell$, to the KL divergence between two probability distributions, with predictions and activities serving as sufficient statistics. In detail, the PC model now has the following variational free energy:

$$\mathcal{E}_t = \sum_{\ell=0}^{L} \mathcal{E}_t^\ell \quad \text{where} \quad \mathcal{E}_t^\ell = D_{KL}[\mathcal{X}^\ell(\mathbf{u}_t^\ell) || \mathcal{X}^\ell(\mathbf{x}_t)], \tag{13}$$

where $\mathcal{X}^\ell(x)$ is a probability distribution defined for every layer with sufficient statistics $x$. If every distribution is a Gaussian, this is equivalent to the energy defined in Equation 4. This generalization allows predictive-coding-based transformers to perform almost as well as standard transformers with the same model complexity [49]. In computational neuroscience, the use of hybrid generative models (with discrete and continuous state spaces) is fairly established: see, e.g., [92] for a first principles account, which can be regarded as a generalization of PC within the setting of well-defined hierarchical generative models with mixed discrete and continuous states.

**Computer Vision.** When processing images, PC has shown promise, particularly when generalizations that utilize convolutional operators are considered to tackle problems in vision such as object recognition and discrimination (using over-complete sparse representations) [96, 97, 98, 99, 100]. Early NGC work used simplifying assumptions like greyscale images and low scene complexity with easily identifiable objects. NGC circuits can process image patches ($P \times P$ pixel grids) to learn low-level patterns such as edges and strokes. NGC was later generalized to natural images through convolutional neural generative coding (Conv-NGC), which integrates standard convolution/deconvolution operations [93]. In Conv-NGC, forward and feedback synapses manipulate feature maps coupled to error units that compute mismatch signals for iterative inference. Notably, Conv-NGC models learn feature representations that implicitly embody image pyramids [101].

**Temporal Data.** PC has achieved a discrete success in temporal sequence modeling. Early work developed hierarchical neural generative modeling for natural video data [102, 99, 103, 104], using top-down information to modulate lower-layer activities for locally invariant representations. This built on PC formulated as a Kalman filter [105, 106] and deep hierarchies of local recurrent networks [107, 108]. One of the earliest NGC formulations addressed time-varying data through the temporal neural coding network (TNCN) [109], later generalized to the parallel temporal neural coding network (P-TNCN) [76]. Both models processed pattern sequences without temporal unfolding, making them local in space and time. They demonstrated zero-shot generalization to unseen sequences, with TNCN adapting to bouncing

| Research Domain | Tasks | References |
|---|---|---|
| Discriminative Learning | Categorization, Graph-based classification | [53, 57, 60] |
| Natural Language Processing | Language modeling, Masked language modeling | [76, 49, 127] |
| Computer Vision | Unsupervised reconstruction, Object classification, Feature extraction | [107, 20, 98, 99, 100, 96, 97, 93] |
| Temporal Modeling | Video sequence modeling, Sequence-based classification | [105, 106, 102, 99, 103] [104, 109, 76, 110] |
| Lifelong Learning | Continual classification, Continual sequence modeling | [75, 76, 111, 67] |
| Control / Robotics | Reinforcement learning, Robotic control, maze navigation | [122, 128, 129, 119, 125, 95] [123, 94, 120, 124, 126] |

Table 1: Overview of PC utilized and developed for machine learning, grouped by application/topic area. For a comprehensive study of the state of the art results in computer vision, we refer to [81].

ball videos of different object quantities and P-TNCN synthesizing frames with varying quantities and qualities without altering synaptic weights. Recently, dynamic PC [110] combined deep learning tools with mixture models, acquiring bio-plausible receptive fields for sequential visual perception tasks.

**Continual Learning.** Another promising aspect of P-TNCNs described above is its ability to conduct continual sequential learning [76], as its generative ability on previously seen sequence modeling tasks did not deteriorate nearly as much as in recurrent networks. Additional effort [75] explored strengthening the memory retention ability of NGC by examining the challenging problem of online cumulative learning, where datasets (or tasks) were presented to the system in the form of a stream and with no indication of when the task was switched. Surprisingly, the NGC-based model was shown to outperform or be competitive with a wide swath of backprop-based approaches that relied on memory-buffers (replay), regularization, and/or auxiliary generative modeling in order to preserve prior knowledge. Stronger performance was seen when another neural circuit drove the lateral recurrent synapses of S-NCN based on competitive learning. This simple task mediating circuit was later improved in follow-up effort [111]. Another later effort developed a stochastic implementation of PC with readable/writable memory (BayesPCN) [67] that was shown to be robust to sample-level forgetting.

**Active Inference and Control.** In computational neuroscience, early formulations of active inference were based on equipping PC with reflexes in order to simulate a variety of behaviors; ranging from handwriting and its observation [112], through oculomotor control [113, 114] to communication [115, 116], inspiring works that merge it with kinematic control [117, 118]. Active NGC (ANGC) [119] and active PC (ActPC) [120] formulate active inference through predictive processing, contrasting with backprop-trained networks. ANGC implements modular agent design with policy and transition dynamics circuits learning from sparse rewards, while ActPC adds actor and prior preference circuits with working memory integration. Both work well on standard RL control and simulated robotics [121, 122, 123]. Parallel frameworks include active PC (APC) [124] and active PC networks (APCNs) [125] for learning part-whole hierarchies from natural images and hierarchical reinforcement learning in multi-room grid-worlds. However, these utilize backprop in portions of their architectures. The COGnitive Neural GENerative system (CogNGen) [95, 94, 126] combines NGC with hyperdimensional memory models, implementing motor cortex, procedural memory, working memory buffers, and episodic memory modules, performing well on maze exploration tasks requiring cross-episode memory.

# 5 Plausibility of Predictive Coding in Neuroscience

A common question in machine learning is: at what point can a particular algorithm be considered to be biologically plausible? This arises from the fact that no computer simulation can fully replicate the intricate workings of the brain in every respect, and, hence, there will invariably be certain nuances that render the simulation implausible in some way. Furthermore, different research agendas consider different properties for differentiating biologically plausible and

non-plausible models. In this section, we will start by addressing the neuroscientific debate and problems with testing PC and then discuss the key properties of PC that satisfy this distinction and which do not.

**Neuroscientific Debate over Predictive Coding.** Despite the fact that PC is a subset of variational Bayesian methodology, there is far less direct experimental support for PC compared to the broader Bayesian brain hypothesis. While there is abundant evidence for the Bayesian brain hypothesis itself [130, 131, 132, 133, 134, 135], direct experimental support for PC across modalities remains more limited [136]. Neuroimaging does support some elements of PC [137, 138, 139], particularly through the framing of mismatch negativity [140, 141]. However, mismatch negativity can also be explained via neural adaptation without invoking hierarchical Bayesian mechanics [136]. Evidence has emerged indicating that PC might be better considered as a cognitive, higher-level explanation rather than a lower-level sensory-based one [136]. Although PC can account for various neurophysiological phenomena such as bistable perception [142], perceptual illusions [143], and properties of V1 neuronal responses [144], fewer studies explicitly test hypotheses inherent to PC. These core hypotheses include expectation-scaled error-signaling neural responses and top-down signaling representing sensory prediction. PC has been criticized as difficult to falsify [145].

The translation of PC from algorithmic specification to biophysical instantiation is often unclear, creating further ambiguity about which implementation is being tested in particular studies. Some PC variants better explain certain neurophysiological data over others [146], making it difficult for neuroscience researchers to test its validity [138]. Some studies suggest that PC is more likely to occur in certain brain regions depending on the sensory modalities being processed [147]. More importantly, there appears to be no evidence of a distinction between predictions and errors at the network level [147, 148]. Nevertheless, despite this unsettled debate surrounding PC and its empirical grounding, advances in neurophysiological methodology offer promise. Higher resolution fMRI, calcium imaging, and other techniques, when paired with anatomical models and simulation, might enable future derivations of finer-grained hypotheses. These advances could also facilitate comparisons across variant PC implementations that share the same underlying theoretical framework, including those with different forms of encoding or representation of prediction errors.

**Error Neurons.** Our understanding of how PC might be implemented at the neuronal level has certainly changed over the last decade. The initial assumption was that the brain would encode two families of neurons/structures: one in charge of propagating predictions and one in charge of propagating errors [20, 53]. As of today, we have no definitive empirical evidence of the existence of single error neurons, though there is ample evidence for laminar-specific segregation of neuronal populations that may communicate predictions and errors, respectively [149, 150]. Recent work has shown that error signals could potentially be computed by the local voltage dynamics in the dendrites [151, 54, 152, 153]. For a more detailed description in the context of a possible biological neural implementation, see [55, 152]. Although existing technologies make it difficult to empirically demonstrate PC at the level of a single neuron, we have stronger evidence of PC at the level of neural populations, or brain regions [138, 136]. PC can account for multiple brain phenomena, such as end-stopping and extra-classical receptive fields effects in V1 [20], bistable perception [154], illusory motion [143], repetition-suppression [155], and attentional modulation of cortical processing [156, 157]. It has also been shown that, when listening to speech, the human brain stores acquired information in a hierarchical fashion, where frontoparietal cortices predict the activities of higher level representations [139].

**Precision Engineering.** Another challenge in the implementation of PC is the management and updating of precision weighting. Typically, the covariance matrix $\Sigma^\ell$ or its inverse, the precision matrix $(\Sigma^\ell)^{-1}$, must be adjusted or computed using matrix inversion [105, 24, 64]. However, implementations of PC in computational neuroscience have solved this problem using standard solutions from precision or covariance component analysis in a biologically plausible way: see [37] (Equation 57) and [158, 56] for details. This is an important aspect of PC, because precision weighting is thought to implement attention in a neuroscience setting. There is a substantial literature on PC in this setting that might usefully inform machine learning implementations [156, 159, 160, 157, 161, 162].

**Synaptic Constraints.** In standard PC networks, synaptic values, which are crucial for neural computation, are generally not subject to strict limitations after an update. This can cause the values to become highly positive or negative, which makes the model less stable overall. This is tackled in NGC, which introduces a constraint to ensure that the Euclidean norm of any row or column in the synaptic matrix does not exceed one, aiding in stability, a practice rooted in early classical sparse coding linear generative models [163]. Beyond this, there have been recent attempts to establish heuristic limits on the magnitude of synaptic values in order to thwart declines in classification performance [164]. A second implausibility of PC models is the frequent alterations in the signs of their synapses (or "sign-flipping"), that can go from negative to positive (and vice versa) during training, an aspect that is pivotal in emulating real cortical functionality. A potential solution to prevent sign-flipping would be to enforce non-negativity constraints on synaptic
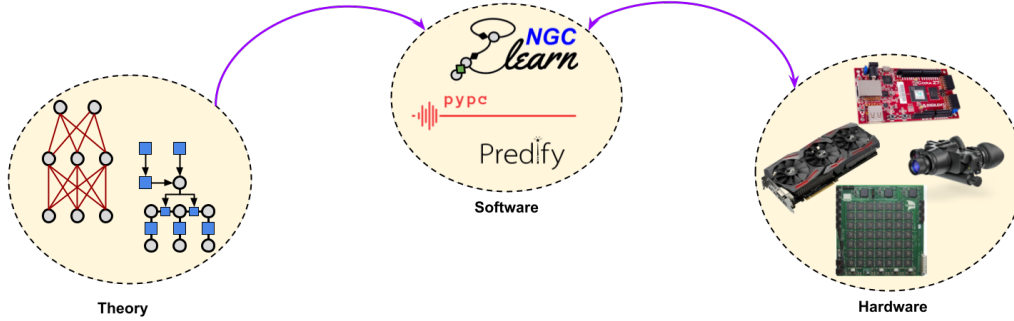
Figure 7: The theory-software-hardware cycle of predictive processing showing bidirectional relationships between theoretical development, simulation software, and hardware platforms.

values and explicitly model groups of excitatory and inhibitory neurons, given that this configuration would apply the necessary positive and negative forces essential for PC inference and learning [165, 166, 167, 168].

# 6 Software Frameworks and Novel Hardware

Software and hardware play crucial roles in enabling innovation and practical implementation of computational intelligence models. While frameworks like PyTorch [169] and TensorFlow [170] have been instrumental for deep learning, PC research faces unique challenges due to sparse software support and evolving hardware opportunities. In terms of software, available frameworks remain relatively sparse, with most research resulting in scattered, paper-specific code that hinders wider adoption and reproducibility. Nevertheless, several libraries aim to democratize PC research: *ngc-learn* [171] is the official library for NGC, is grounded in neuronal cable theory and supports arbitrary PC model and neuromorphic system construction as well as general neural circuit simulation. *PCX* [81] is a deep learning oriented library, is built on JAX/Equinox and enables plug-and-play deep learning with PC updates and demonstrates state-of-the-art results across frameworks. Other libraries include *pypc* [172] for hierarchical PC models, *predify* [173] for converting deep architectures to PC-like systems, and *pyhgf* [174] for stochastic and Bayesian PC formulations, and the recently proposed jPC [87] for fully parallel deep learning architectures.

## 6.1 Novel Hardware Implementation

Hardware limitations historically shape research directions [175]. While GPUs and TPUs dominate current training, emerging technologies like memristors, spintronics, and optics could revolutionize the field [176, 177, 178]. PC is particularly well-suited for alternative hardware as an energy-based model trainable through equilibrium propagation [15, 179]. Its iterative inference with layer-wise parallel computation, Hebbian adaptation, and local energy optimization enables significant parallelization, reducing communication bottlenecks inherent in backpropagation. This alignment with neuromorphic hardware creates opportunities for energy-efficient learning [180], with dynamical memristors offering promising potential for brain-inspired systems [181]. Beyond conventional hardware, PC could enable "intelligence-in-a-dish" technology through adaptive computation in biophysical mediums [182, 183, 184]. Given the association of PC with cortical computation [144, 149, 185], it is inherently compatible with organoid intelligence [186]. Early spike-level implementations [187, 165] show promise for dynamic predictive processing. Organoid growth and decay could facilitate model selection [188], potentially improving generalization through natural evolution of neural structures.

**Spiking Predictive Coding.** A key challenge involves generalizing PC to spike-level processing. Biological neurons communicate through sparse action potentials creating information-rich spike trains [189], inspiring spiking neural networks that encode information through precise timing [190]. This enables energy-efficient neuromorphic hardware [191, 192, 193], unlike energy-intensive GPU-based networks. However, most PC formulations fail to account for sparse, discrete communication, reflecting the nature of PC as describing ensemble rather than individual neuron dynamics [131]. The spiking neural coding framework attempts PC implementation in spiking networks using spike-induced synaptic adjustments [194, 165]. Unlike standard spike-timing-dependent plasticity requiring hand-crafted layers [195], this framework uses flexible synaptic conductance models, membrane potentials, and trace mechanisms. Each layer estimates the trace activity of the following layer, passing discrepancies through evolved error feedback synapses (Fig. 8). The framework adapts to various spiking functions from leaky integrators to Hodgkin-Huxley models
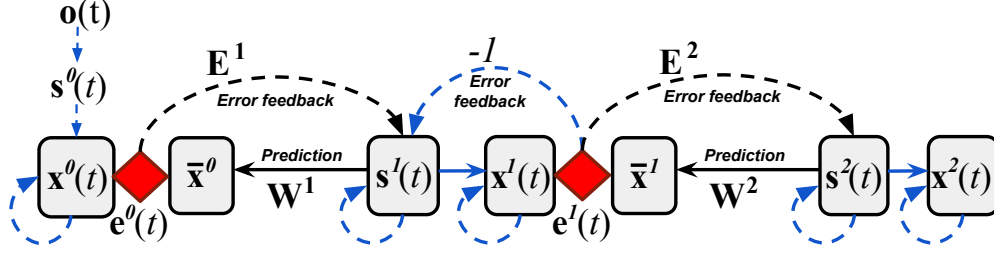
Figure 8: Spiking neural coding formulation where $\mathbf{s}^{\ell}(t)$ represents spike nodes and $\mathbf{x}^{\ell}(t)$ are trace approximations. Solid arrows show forward synapses, dashed arrows show error transmission.

[196]. Recent implementations [197, 198, 152] operate in discriminative contexts but do not address weight transport, sharing feedback with forward synapses. For comprehensive coverage, see [199].

# 7 Looking Forward: Important Directions for Predictive Coding Research

As we have discussed in Sections 3 and 4, the performance of PC models can be quite below that of modern-day deep neural networks trained with backprop. The future of PC in machine learning strongly depends on our ability to address and fill this gap. Specifically, among the main directions for future work related to PC should be understanding the fundamental reasons behind this performance mismatch and using this acquired insight to develop and design novel PC schemas, mathematical frameworks, and computational models that work in large-scale settings, where current deep learning models excel.

**Efficiency.** The first drawback of PC is its efficiency. This is a consequence of its underlying iterative inference process, which generally needs to be run until convergence. In practice, it is common to let a PC model run for a fixed number of iterations $T$, but this number must be large in order to reach high performance, and deeper networks require more iterations in order to perform well. To this end, it would be useful to derive different optimization techniques and methodologies that carry out the minimization of the variational free energy. Such techniques could come in the form of (faster) variations of the expectation-maximization (EM) algorithm (equivalent to gradient descent in the standard case). An example of such a scheme is the newly proposed *incremental* PC [57], based on incremental EM. The convergence of this method has also been proven using methods from dynamical systems, under the assumption of sufficiently small parameters [200]. Future work could investigate whether there are better alternatives or whether the optimal update rule can be learned with respect to specific tasks and datasets (drawing inspiration from the recent developments in deep meta-learning [201, 202, 203]). More advanced amortized inference algorithms and mechanisms might prove to be another critical ingredient for reducing the computational burden of the PC inference process itself, as, historically, the iterative inference of classical sparse coding linear generative models [163] was dramatically sped up when introducing a second recognition neural model [204]. It might be that the answer may lie in some variations of the EM algorithm that are lesser known to the machine learning community, such as *dynamical* expectation maximization (DEM) [205, 43, 43], or in alternative, more efficient implementations of precision-weighted prediction errors, as in [110]. Other possibilities may lie in methods based on different message passing frameworks, such as belief propagation and variational message passing over factor graphs.

**Optimization Tricks and Heuristics.** Future research will also need to focus on the study of optimization techniques that have proven useful and invaluable for variational inference, such as those that promote the inclusion of precision weighting parameters into the picture [206]. In the same spirit, recent progress on natural-gradient methods for gradient-based variational inference—ranging from stochastic NGVI with non-asymptotic guarantees, through Kronecker-factored curvature approximations, to the variational predictive natural gradient—offers curvature-aware updates that can be interpreted as learning layer-wise precisions online, and therefore provides an immediately transferable optimisation tool-set for scaling precision-weighted predictive coding [207, 208, 209]. Despite interesting developments obtained over the last several years, and despite the fact that precisions are of primary importance in simulations used in the neurosciences, these techniques have only been tested in small and medium scale settings [210, 211, 212]. More generally, the field of deep learning has enjoyed immense benefits from simple optimization tricks developed throughout the last decade, such as dropout [213], batch normalization [214], adaptive learning rates such as Adam [215] and RMSprop [216], and the introduction of the ReLU activation [217] and its variants. Without these techniques, emerging from the joint effort of thousands of researchers, the training of extremely overparametrized neural models would not lead to the results that we observe today. Interestingly, backpropagation itself exhibited many flaws and limitations at

14

the beginning that precluded it scaling to high-dimensional data spaces: in the late 90s, kernel learning methods were dominant and generally found to be more effective than artificial neural networks, which were generally disregarded due the problem of vanishing gradients [218] as well as their strong requirement for heavy computational power. Most of these problems have now been addressed in over thirty years of research, again through a collective effort of thousands of people. Given the history of deep learning, one might ask: what are the dropout, batch norm, and Adam optimizer equivalents for PC? Addressing such a question would be of vital importance if the goal is to scale the applicability of PC, and will hopefully become a more prominent topic of interest of more research efforts going forward.

**Stochastic Generative Models and Sampling.** A different but important, direction, is that related to generative models. We have been discussing the importance of putting statistical models of joint distributions back into the research agenda of generative AI; however, most of the research efforts that go beyond small-scale tasks focus on supervised learning [57, 49] (with small notable exceptions, such as the unsupervised neural circuit of [93]). To do so, we need to develop models that are capable of sampling data points from a well-computed posterior distribution, using simulation methods such those based on Langevin dynamics [219]. This would benefit different subfields of Bayesian inference, such as out-of-distribution (OOD) detection, uncertainty minimization, and data reconstruction. In fact, PC is particularly suitable to OOD detection, thanks to a measure of surprise that is always readily available to the model, that is, with respect to its variational free energy.

Another probabilistic generative modeling research direction for PC is to incorporate the uncertainty about model (i.e., synaptic) parameters into the variational bound on marginal likelihood. That is, one can equip synaptic parameters with probability distributions (as opposed to using point estimates) [220]. This move would take PC closer to its Bayesian roots: current machine learning implementations of PC do not treat the model parameters as random variables and can therefore be regarded as an expectation maximization (EM) procedure, where the M-step ignores uncertainty about the parameters [221]. The benefit of treating parameters as random variables is that one can evaluate the model evidence required for structure learning [222, 223, 224, 225], i.e., one could place prediction errors on the synaptic weights and evaluate the resulting variational free energy (or marginal likelihood) for structure learning or, in a biological setting, morphogenesis [226, 227]. This research direction takes PC into the world of generalized filtering and dynamic causal modeling, in which parameters can be regarded as latent states that change very slowly [228, 229]. In the case of dynamic causal modeling (i.e., the variational inversion of state-space models under Gaussian assumptions), structure learning can be rendered particularly efficient through particular instances of Bayesian model selection, which provides a principled basis for comparing alternative generative hypotheses [230, 231, 232]. In predictive coding (PC), the same logic appears as *Bayesian model reduction*, where synaptic weights that receive little posterior evidence are pruned from a full (parent) network [233]. Looking ahead, a natural extension is to replace these fixed-dimension reductions with **Bayesian non-parametric** (BNP) priors—such as Dirichlet or Indian-buffet processes—that equip each PC layer with an *unbounded* supply of latent causes or synaptic motifs [231, 234, 235]. Thus in BNP–PC hybrid, precision-weighted prediction errors might simultaneously (i) adjust beliefs about existing components and (ii) signal when residual error warrants "breaking a new stick" (creating a component) or dropping a redundant one. Thus, BNP priors transform PC into a self-organising system that learns both its parameters *and* its effective dimensionality online, unifying synaptic growth and pruning within a single, neurally plausible free-energy minimisation scheme.

**Crafting Control Systems.** PC brings with it the promise of learning a powerful generative model that is continuously and iteratively refined as more sensory samples are gathered over time. This has led some early work to consider such a process as the basis for world models that drive modular, brain-motivated cognitive models, capable of combining perception and action in the context of playing video games [96, 129] and robotic control tasks [121, 119, 120], as well as large-scale cognitive architectures [236, 94, 126]. This has important implications beyond machine learning, particularly for the domains of cognitive science and cognitive neuroscience where a key pathway is to craft computational theories of mind and to examine their fit to human subject data on controlled psychological tasks [237] as well as their ability to generalize to cognitive functionality. Hence, a promising future direction is to design modular, increasingly complex systems made up of PC circuitry, that rely on single fundamental PC circuits and low-level dynamics.

There are important lessons that come from crafting systems of PC circuitry that can inform the fundamentals of PC itself. For example, in the PC-centric cognitive architecture of [94], it was found that the synergy of PC with another important neural model, i.e., vector symbolic memory [238, 239, 240], facilitated effective complex auto-associative and hetero-associative memory operations that led to rapid convergence on complex maze navigation tasks and, furthermore, facilitated the design of a novel PC circuit that leveraged its neural dynamics to learn a policy for internally manipulating a flexible recurrent memory system [126].

Using PC as a fundamental neural building block for control systems is still in its earliest stages, and although more effort will be needed for it to become viable for building robust computational theories of mind, it presents a worthwhile long-term direction for crafting simulated natural intelligence. This promise is made even stronger when considering that

deep backprop-based networks are beginning to find use in the field of cognitive science [241]. Making developments along this direction would also crucially benefit the research in active inference [242, 92, 19, 18], given that agent systems that engage in epistemic foraging often center around the use of dynamic generative models and can even be viewed as simple control models.

## 8 Conclusion and Outlook

Generative models will most likely play a large role in the future of artificial intelligence. Despite this, current research seems to mostly focus on a limited class of models, overlooking possible alternatives. On the one hand, this is justified by the exceptional results obtained by deep, backprop-trained artificial neural networks; on the other hand, more principled methods, whose formulations are grounded in statistics and information theory, may deserve attention. In this work, we have focused on inverting a specific class of continuous-state generative models, namely, predictive coding. More precisely, we have: (1) presented a review of predictive coding schemes in machine learning, highlighting previous art that has led us to the point that we are at today, (2) summarized important open questions in the field that need to be answered in order to unlock and utilize the full potential of predictive coding, and (3) identified possible applications and future research directions. As of today, we are seeing some of the potential of PC, but we are still far away from large-scale applications that may call for a significant investment in PC research. Already in the past years, we have seen impressive growth in the field of predictive coding in machine learning: it was only in 2017 that Whittington and Bogacz demonstrated how to train a small predictive coding classifier on MNIST [53]. In addition, we have already seen how integrating a few ideas from predictive coding can lead to interesting and more powerful backprop-based deep models [243, 244, 58, 245, 246, 127].

One of the main goals of this survey is to encourage researchers to build on the results of decades of prior effort and focus on the challenges offered by predictive coding: no success story can come without an active community behind it, as evidenced by backpropagation-centric deep learning. We remark that a promising methodology is only as good as the efforts that are made to advance it, both empirically and theoretically. Community effort will be needed to advance predictive coding from both the software and the hardware standpoint; particularly, to develop computational schemes that exploit the advantages on offer, such as its parallelism and sparse, local, and potentially energy-efficient computations. Although research progress in predictive coding has been consistent over the past decades, we may be only starting to realize the benefits afforded to artificial intelligence by reverse engineering the cortex and other biological structures.

## Acknowledgments

## References

[1] Rombach, R, Blattmann, A, Lorenz, D, Esser, P, and Ommer, B. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022.

[2] Saharia, C, Chan, W, Saxena, S, Li, L, Whang, J, Denton, EL, Ghasemipour, K, Gontijo Lopes, R, Karagol Ayan, B, Salimans, T, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022.

[3] Ho, J, Chan, W, Saharia, C, Whang, J, Gao, R, Gritsenko, A, Kingma, DP, Poole, B, Norouzi, M, Fleet, DJ, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.

[4] Ramesh, A, Dhariwal, P, Nichol, A, Chu, C, and Chen, M. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.

[5] (FAIR)†, MFARDT, Bakhtin, A, Brown, N, Dinan, E, Farina, G, Flaherty, C, Fried, D, Goff, A, Gray, J, Hu, H, et al. Human-level play in the game of diplomacy by combining language models with strategic reasoning. *Science*, 378(6624):1067–1074, 2022.

[6] Perolat, J, De Vylder, B, Hennes, D, Tarassov, E, Strub, F, de Boer, V, Muller, P, Connor, JT, Burch, N, Anthony, T, et al. Mastering the game of stratego with model-free multiagent reinforcement learning. *Science*, 378(6623):990–996, 2022.

[7] Touvron, H, Martin, L, Stone, K, Albert, P, Almahairi, A, Babaei, Y, Bashlykov, N, Batra, S, Bhargava, P, Bhosale, S, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

[8] Chen, M, Tworek, J, Jun, H, Yuan, Q, Pinto, HPdO, Kaplan, J, Edwards, H, Burda, Y, Joseph, N, Brockman, G, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.

[9] Brown, TB, Mann, B, Ryder, N, Subbiah, M, Kaplan, J, Dhariwal, P, Neelakantan, A, Shyam, P, Sastry, G, Askell, A, et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 2020.

[10] Ciregan, D, Meier, U, and Schmidhuber, J. Multi-column deep neural networks for image classification. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 3642–3649. IEEE, 2012.

[11] Krizhevsky, A, Sutskever, I, and Hinton, GE. ImageNet classification with deep convolutional neural networks. In *26th Annual Conference on Neural Information Processing Systems (NIPS) 2012*, 2012.

[12] Linnainmaa, S. *The representation of the cumulative rounding error of an algorithm as a Taylor expansion of the local rounding errors*. PhD thesis, Master's Thesis (in Finnish), Univ. Helsinki, 1970.

[13] Rumelhart, DE, Hinton, GE, and Williams, RJ. Learning representations by back-propagating errors. *Nature*, 323(6088):533–536, 1986.

[14] Bengio, Y, Lee, DH, Bornschein, J, Mesnard, T, and Lin, Z. Towards biologically plausible deep learning. *arXiv preprint arXiv:1502.04156*, 2015.

[15] Scellier, B and Bengio, Y. Equilibrium propagation: Bridging the gap between energy-based models and backpropagation. *Frontiers in Computational Neuroscience*, 11:24, 2017.

[16] Hinton, G. The forward-forward algorithm: Some preliminary investigations. *arXiv preprint arXiv:2212.13345*, 2022.

[17] Zador, A, Richards, B, Ölveczky, B, Escola, S, Bengio, Y, Boahen, K, Botvinick, M, Chklovskii, D, Churchland, A, Clopath, C, et al. Toward next-generation artificial intelligence: Catalyzing the NeuroAI revolution. *arXiv preprint arXiv:2210.08340*, 2022.

[18] Friston, KJ, Ramstead, MJD, Kiefer, AB, Tschantz, A, Buckley, CL, Albarracin, M, Pitliya, RJ, Heins, C, Klein, B, Millidge, B, et al. Designing ecosystems of intelligence from first principles. *arXiv preprint arXiv:2212.01354*, 2022.

[19] Da Costa, L, Lanillos, P, Sajid, N, Friston, K, and Khan, S. How active inference could help revolutionise robotics. *Entropy*, 24(3):361, 2022.

[20] Rao, RPN and Ballard, DH. Predictive coding in the visual cortex: A functional interpretation of some extra-classical receptive-field effects. *Nature Neuroscience*, 2(1):79–87, 1999.

[21] Friston, K and Kiebel, S. Predictive coding under the free-energy principle. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 364(1521):1211–1221, 2009.

[22] Millidge, B, Salvatori, T, Song, Y, Bogacz, R, and Lukasiewicz, T. Predictive coding: Towards a future of deep learning beyond backpropagation? In *Proceedings of the 31st International Joint Conference on Artificial Intelligence and the 25th European Conference on Artificial Intelligence, IJCAI-ECAI 2022, Survey Track, Vienna, Austria, July 23–29, 2022*, pages 5538–5545. IJCAI/AAAI Press, July 2022.

[23] Hebb, D. *The Organization of Behavior*. Wiley, New York, 1949.

[24] Ororbia, A and Kifer, D. The neural coding framework for learning generative models. *Nature Communications*, 13(1):2064, 2022.

[25] Salvatori, T, Pinchetti, L, Millidge, B, Song, Y, Bao, T, Bogacz, R, and Lukasiewicz, T. Learning on arbitrary graph topologies via predictive coding. *Advances in Neural Information Processing Systems*, 2022.

[26] Avena-Koenigsberger, A, Misic, B, and Sporns, O. Communication dynamics in complex brain networks. *Nature Reviews Neuroscience*, 19(1):17–33, 2018.

[27] Song, Y, Millidge, BG, Salvatori, T, Lukasiewicz, T, Xu, Z, and Bogacz, R. Inferring neural activity before plasticity: A foundation for learning beyond backpropagation. *Nature Neuroscience*, 2023.

[28] Alonso, N, Millidge, B, Krichmar, J, and Neftci, EO. A theoretical framework for inference learning. *Advances in Neural Information Processing Systems*, 35:37335–37348, 2022.

[29] Innocenti, F, Singh, R, and Buckley, CL. Understanding predictive coding as an adaptive trust-region method. *arXiv preprint arXiv:2305.18188*, 2023.

[30] Kingma, DP and Welling, M. Auto-encoding variational Bayes. *arXiv preprint arXiv:1312.6114*, 2013.

[31] Domingos, P. The role of Occam's razor in knowledge discovery. *Data Mining and Knowledge Discovery*, 3:409–425, 1999.

[32] Cutler, CC. Differential PCM, U.S. Patent 2 605 361, July 29, 1952.

[33] O'Neal, J. Entropy coding in speech and television differential PCM systems (corresp.). *IEEE Transactions on Information Theory*, 17(6):758–761, 1971.

[34] Elias, P. Predictive coding—I. *IRE Transactions on Information Theory*, 1(1):16–24, 1955.

[35] Srinivasan, MV, Laughlin, SB, and Dubs, A. Predictive coding: A fresh view of inhibition in the retina. *Proceedings of the Royal Society of London. Series B. Biological Sciences*, 216(1205):427–459, 1982.

[36] Friston, K. Learning and inference in the brain. *Neural Networks*, 16(9):1325–1352, 2003.

[37] Friston, K. Hierarchical models in the brain. *PLoS Computational Biology*, 2008.

[38] Friston, K. The free-energy principle: A unified brain theory? *Nature Reviews Neuroscience*, 11(2):127–138, 2010.

[39] Winn, J, Bishop, CM, and Jaakkola, T. Variational message passing. *Journal of Machine Learning Research*, 6(4), 2005.

[40] Elias, P. Predictive coding—II. *IRE Transactions on Information Theory*, 1(1):16–24, 1955.

[41] Friston, K. A theory of cortical responses. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 360(1456), 2005.

[42] Spratling, MW. A review of predictive coding algorithms. *Brain and Cognition*, 112:92–97, 2017.

[43] Friston, KJ, Trujillo-Barreto, N, and Daunizeau, J. DEM: A variational treatment of dynamic systems. *Neuroimage*, 41(3):849–885, 2008.

[44] Millidge, B, Tschantz, A, and Buckley, CL. Predictive coding approximates backprop along arbitrary computation graphs. *arXiv:2006.04182*, 2020.

[45] Vafaii, H, Galor, D, and Yates, J. Poisson variational autoencoder. *Advances in Neural Information Processing Systems*, 37:44871–44906, 2024.

[46] Vafaii, H, Galor, D, and Yates, JL. Brain-like variational inference. *ArXiv*, pages arXiv–2410, 2025.

[47] Zahid, U, Guo, Q, and Fountas, Z. Sample as you infer: Predictive coding with Langevin dynamics. *arXiv preprint arXiv:2311.13664*, 2023.

[48] Tscshantz, A, Millidge, B, Seth, AK, and Buckley, CL. Hybrid predictive coding: Inferring, fast and slow. *PLoS computational biology*, 19(8):e1011280, 2023.

[49] Pinchetti, L, Salvatori, T, Millidge, B, Song, Y, Yordanov, Y, and Lukasiewicz, T. Predictive coding beyond Gaussian assumptions. *Advances in Neural Information Processing Systems*, 2022.

[50] Salvatori, T, Pinchetti, L, M'Charrak, A, Millidge, B, and Lukasiewicz, T. Predictive coding beyond correlations. *arXiv preprint arXiv:2306.15479*, 2023.

[51] Sennesh, E, Wu, H, and Salvatori, T. Divide-and-conquer predictive coding: A structured Bayesian inference algorithm. *arXiv preprint arXiv:2408.05834*, 2024.

[52] Oliviers, G, Bogacz, R, and Meulemans, A. Learning probability distributions of sensory inputs with Monte Carlo predictive coding. *PLoS Computational Biology*, 20(10):e1012532, 2024.

[53] Whittington, JCR and Bogacz, R. An approximation of the error backpropagation algorithm in a predictive coding network with local Hebbian synaptic plasticity. *Neural Computation*, 29(5), 2017.

[54] Sacramento, J, Costa, RP, Bengio, Y, and Senn, W. Dendritic cortical microcircuits approximate the backpropagation algorithm. In *Advances in Neural Information Processing Systems*, pages 8721–8732, 2018.

[55] Whittington, JCR and Bogacz, R. Theories of error back-propagation in the brain. *Trends in Cognitive Sciences*, 2019.

[56] Millidge, B, Seth, A, and Buckley, CL. Predictive coding: A theoretical and experimental review. *arXiv:2107.12979*, 2021.

[57] Salvatori, T, Song, Y, Millidge, B, Xu, Z, Sha, L, Emde, C, Bogacz, R, and Lukasiewicz, T. Incremental predictive coding: A parallel and fully automatic learning algorithm. *arXiv preprint arXiv:2212.00720*, 2022.

[58] Han, K, Wen, H, Zhang, Y, Fu, D, Culurciello, E, and Liu, Z. Deep predictive coding network with local recurrent processing for object recognition. *Advances in Neural Information Processing Systems*, 31, 2018.

[59] Millidge, B, Song, Y, Salvatori, T, Lukasiewicz, T, and Bogacz, R. A theoretical framework for inference and learning in predictive coding networks. In *Proceedings of the 11th International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, 1–5 May 2023*. OpenReview.net, 2023.

[60] Byiringiro, B, Salvatori, T, and Lukasiewicz, T. Robust graph representation learning via predictive coding. *arXiv preprint arXiv:2212.04656*, 2022.

[61] Papadimitriou, CH, Vempala, SS, Mitropolsky, D, Collins, M, and Maass, W. Brain computation by assemblies of neurons. *Proceedings of the National Academy of Sciences*, 2020.

[62] Rao, RPN. An optimal estimation approach to visual perception and learning. *Vision Research*, 39(11):1963–1989, 1999.

[63] Salvatori, T, Song, Y, Hong, Y, Sha, L, Frieder, S, Xu, Z, Bogacz, R, and Lukasiewicz, T. Associative memories via predictive coding. In *Advances in Neural Information Processing Systems*, volume 34, 2021.

[64] Tang, M, Salvatori, T, Millidge, B, Song, Y, Lukasiewicz, T, and Bogacz, R. Recurrent predictive coding models for associative memory employing covariance learning. *PLoS Computational Biology*, 19(4):e1010719, 2023.

[65] Millidge, B, Salvatori, T, Song, Y, Lukasiewicz, T, and Bogacz, R. Universal Hopfield networks: A general framework for single-shot associative memory models. In *Proceedings of the 39th International Conference on Machine Learning, ICML 2022, Baltimore, Maryland, USA, 17-23 July 2022*, Proceedings of Machine Learning Research, pages 15561–15583. PMLR, 2022.

[66] Ramsauer, H, Schäfl, B, Lehner, J, Seidl, P, Widrich, M, Gruber, L, Holzleitner, M, Adler, T, Kreil, D, Kopp, MK, Klambauer, G, Brandstetter, J, and Hochreiter, S. Hopfield networks is all you need. In *International Conference on Learning Representations*, 2021.

[67] Yoo, J and Wood, F. BayesPCN: A continually learnable predictive coding associative memory. *Advances in Neural Information Processing Systems*, 35:29903–29914, 2022.

[68] Annabi, L, Pitti, A, and Quoy, M. On the relationship between variational inference and auto-associative memory. *Advances in Neural Information Processing Systems*, 2022.

[69] Krotov, D and Hopfield, JJ. Large associative memory problem in neurobiology and machine learning. In *International Conference on Learning Representations*, 2021.

[70] Millidge, B, Song, Y, Salvatori, T, Lukasiewicz, T, and Bogacz, R. Backpropagation at the infinitesimal inference limit of energy-based models: Unifying predictive coding, equilibrium propagation, and contrastive Hebbian learning. In *Proceedings of the 11th International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, 1–5 May 2023*, 2023.

[71] Rosenbaum, R. On the relationship between predictive coding and backpropagation. *PLoS One*, 17(3):e0266102, 2022.

[72] Song, Y, Lukasiewicz, T, Xu, Z, and Bogacz, R. Can the brain do backpropagation? — Exact implementation of backpropagation in predictive coding networks. In *Advances in Neural Information Processing Systems*, volume 33, 2020.

[73] Salvatori, T, Song, Y, Xu, Z, Lukasiewicz, T, and Bogacz, R. Reverse differentiation via predictive coding. In *Proceedings of the 36th AAAI Conference on Artificial Intelligence*. AAAI Press, 2022.

[74] Zahid, U, Guo, Q, and Fountas, Z. Predictive coding as a neuromorphic alternative to backpropagation: A critical evaluation. *Neural Computation*, 35(12):1881–1909, 2023.

[75] Ororbia, A, Mali, A, Giles, CL, and Kifer, D. Lifelong neural predictive coding: Learning cumulatively online without forgetting. *Advances in Neural Information Processing Systems*, 35:5867–5881, 2022.

[76] Ororbia, A, Mali, A, Giles, CL, and Kifer, D. Continual learning of recurrent neural networks by locally aligning distributed representations. *IEEE Transactions on Neural Networks and Learning Systems*, 31(10):4267–4278, 2020.

[77] Barrett, DG and Dherin, B. Implicit gradient regularization. *arXiv preprint arXiv:2009.11162*, 2020.

[78] Toulis, P, Airoldi, E, and Rennie, J. Statistical analysis of stochastic gradient methods for generalized linear models. In *International Conference on Machine Learning*, pages 667–675. PMLR, 2014.

[79] Choksi, B, Mozafari, M, Biggs O'May, C, Ador, B, Alamia, A, and VanRullen, R. Predify: Augmenting deep neural networks with brain-inspired predictive coding dynamics. *Advances in Neural Information Processing Systems*, 34:14069–14083, 2021.

[80] Mali, A, Salvatori, T, and Ororbia, A. Tight stability, convergence, and robustness bounds for predictive coding networks. *arXiv preprint arXiv:2410.04708*, 2024.

[81] Pinchetti, L, Qi, C, Lokshyn, O, Olivers, G, Emde, C, Tang, M, M'Charrak, A, Frieder, S, Menzat, B, Bogacz, R, et al. Benchmarking predictive coding networks–made simple. *arXiv preprint arXiv:2407.01163*, 2024.

[82] Koudahl, M, Salvatori, T, Da Costa, L, Beck, J, Tschantz, A, Buckley, CL, and Linander, H. Bethe predictive coding. *International Active Inference Workshop IWAI*, 2025.

[83] Innocenti, F, Achour, EM, and Buckley, CL. $\mu$pc: Scaling predictive coding to 100+ layer networks. *arXiv preprint arXiv:2505.13124*, 2025.

[84] Qi, C, Lukasiewicz, T, and Salvatori, T. Training deep predictive coding networks. In *New Frontiers in Associative Memories*, 2025.

[85] Goemaere, C, Oliviers, G, Bogacz, R, and Demeester, T. Error optimization: Overcoming exponential signal decay in deep predictive coding networks. *arXiv preprint arXiv:2505.20137*, 2025.

[86] Devlin, J, Chang, MW, Lee, K, and Toutanova, K. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics*. Association for Computational Linguistics, 2019.

[87] Innocenti, F, Kinghorn, P, Yun-Farmbrough, W, Varona, MDL, Singh, R, and Buckley, CL. JPC: Flexible inference for predictive coding networks in JAX. *arXiv preprint arXiv:2412.03676*, 2024.

[88] Hodgkin, AL and Rushton, WAH. The electrical constants of a crustacean nerve fibre. *Proceedings of the Royal Society of London. Series B-Biological Sciences*, 133(873):444–479, 1946.

[89] Rall, W. Theory of physiological properties of dendrites. *Annals of the New York Academy of Sciences*, 96(4):1071–1092, 1962.

[90] Ermentrout, B and Terman, DH. *Mathematical Foundations of Neuroscience*, volume 35. Springer, 2010.

[91] Vaswani, A, Shazeer, N, Parmar, N, Uszkoreit, J, Jones, L, Gomez, AN, Kaiser, L, and Polosukhin, I. Attention is all you need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems*, 2017.

[92] Friston, KJ, Parr, T, and de Vries, B. The graphical brain: Belief propagation and active inference. *Network Neuroscience*, 1(4):381–414, 2017.

[93] Ororbia, A and Mali, A. Convolutional neural generative coding: Scaling predictive coding to natural images. *Proceedings of the 45th Annual Conference of the Cognitive Science Society*, 2023.

[94] Ororbia, A and Kelly, MA. Cogngen: Building the kernel for a hyperdimensional predictive processing cognitive architecture. In *Proceedings of the Annual Meeting of the Cognitive Science Society*, volume 44, 2022.

[95] Ororbia, A and Kelly, MA. Towards a predictive processing implementation of the common model of cognition. *arXiv preprint arXiv:2105.07308*, 2021.

[96] Principe, JC and Chalasani, R. Cognitive architectures for sensory processing. *Proceedings of the IEEE*, 102(4):514–525, 2014.

[97] Chalasani, R and Principe, JC. Temporal context in object recognition. In *2012 IEEE International Workshop on Machine Learning for Signal Processing*, pages 1–6. IEEE, 2012.

[98] Sledge, IJ and Principe, JC. Faster convergence in deep-predictive-coding networks to learn deeper representations. *IEEE Transactions on Neural Networks and Learning Systems*, 2021.

[99] Chalasani, R and Principe, JC. Context dependent encoding using convolutional dynamic networks. *IEEE Transactions on Neural Networks and Learning Systems*, 26(9):1992–2004, 2014.

[100] Spratling, MW. A hierarchical predictive coding model of object recognition in natural images. *Cognitive Computation*, 9(2):151–167, 2017.

[101] Adelson, EH, Anderson, CH, Bergen, JR, Burt, PJ, and Ogden, JM. Pyramid methods in image processing. *RCA Engineer*, 29(6):33–41, 1984.

[102] Chalasani, R and Principe, JC. Deep predictive coding networks. *arXiv preprint arXiv:1301.3541*, 2013.

[103] Santana, E, Cinar, GT, and Principe, JC. Parallel flow in deep predictive coding networks. In *2015 International Joint Conference on Neural Networks (IJCNN)*, pages 1–5. IEEE, 2015.

[104] Chen, B, Dang, L, Gu, Y, Zheng, N, and Príncipe, JC. Minimum error entropy Kalman filter. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 51(9):5819–5829, 2019.

[105] Rao, RPN and Ballard, DH. Dynamic model of visual recognition predicts neural response properties in the visual cortex. *Neural Computation*, 9(4):721–763, 1997.

[106] Erdogmus, D, Sanchez, JC, and Principe, JC. Modified Kalman filter based method for training state-recurrent multilayer perceptrons. In *Proceedings of the 12th IEEE Workshop on Neural Networks for Signal Processing*, pages 219–228. IEEE, 2002.

[107] Schmidhuber, J. *Neural sequence chunkers*. Inst. für Informatik, 1991.

[108] Schmidhuber, J. Learning complex, extended sequences using the principle of history compression. *Neural Computation*, 4(2):234–242, 1992.

[109] Ororbia, A, Haffner, P, Reitter, D, and Giles, CL. Learning to adapt by minimizing discrepancy. *arXiv:1711.11542*, 2017.

[110] Jiang, LP and Rao, RPN. Dynamic predictive coding: A new model of hierarchical sequence learning and prediction in the cortex. *bioRxiv*, pages 2022–06, 2022.

[111] Ororbia, A. Continual competitive memory: A neural system for online task-free lifelong learning. *arXiv preprint arXiv:2106.13300*, 2021.

[112] Friston, K, Mattout, J, and Kilner, J. Action understanding and active inference. *Biological Cybernetics*, 104:137–160, 2011.

[113] Friston, K, Adams, RA, Perrinet, L, and Breakspear, M. Perceptions as hypotheses: Saccades as experiments. *Frontiers in Psychology*, 3:151, 2012.

[114] Perrinet, LU, Adams, RA, and Friston, KJ. Active inference, eye movements and oculomotor delays. *Biological Cybernetics*, 108:777–801, 2014.

[115] Kiebel, SJ, Daunizeau, J, and Friston, KJ. Perception and hierarchical dynamics. *Frontiers in Neuroinformatics*, 3:569, 2009.

[116] Isomura, T, Parr, T, and Friston, K. Bayesian filtering with multiple internal models: Toward a theory of social intelligence. *Neural Computation*, 31(12):2390–2431, 2019.

[117] Priorelli, M, Pezzulo, G, and Stoianov, IP. Deep kinematic inference affords efficient and scalable control of bodily movements. *Proceedings of the National Academy of Sciences*, 120(51):e2309058120, 2023.

[118] Priorelli, M and Stoianov, IP. Efficient motor learning through action-perception cycles in deep kinematic inference. In *International Workshop on Active Inference*, pages 59–70. Springer, 2023.

[119] Ororbia, AG and Mali, A. Backprop-free reinforcement learning with active neural generative coding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 29–37, 2022.

[120] Ororbia, A and Mali, A. Active predictive coding: Brain-inspired reinforcement learning for sparse reward robotic control problems. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3015–3021. IEEE, 2023.

[121] Tani, J. *Exploring robotic minds: Actions, symbols, and consciousness as self-organizing dynamic phenomena*. Oxford University Press, 2016.

[122] Lanillos, P and Cheng, G. Adaptive robot body learning and estimation through predictive coding. In *Proc. of IROS*, 2018.

[123] Oliver, G, Lanillos, P, and Cheng, G. An empirical study of active inference on a humanoid robot. *IEEE Transactions on Cognitive and Developmental Systems*, 2021.

[124] Rao, RPN, Gklezakos, DC, and Sathish, V. Active predictive coding: A unified neural framework for learning hierarchical world models for perception and planning. *arXiv preprint arXiv:2210.13461*, 2022.

[125] Gklezakos, DC and Rao, RPN. Active predictive coding networks: A neural solution to the problem of learning reference frames and part-whole hierarchies. *arXiv preprint arXiv:2201.08813*, 2022.

[126] Ororbia, AG and Kelly, MA. Maze learning using a hyperdimensional predictive processing cognitive architecture. In *Artificial General Intelligence: 15th International Conference, AGI 2022, Seattle, WA, USA, August 19–22, 2022, Proceedings*, pages 321–331. Springer, 2023.

[127] Araujo, V, Villa, A, Mendoza, M, Moens, MF, and Soto, A. Augmenting bert-style models with predictive coding to improve discourse-level representations. *arXiv preprint arXiv:2109.04602*, 2021.

[128] Muhammad, W and Spratling, MW. A neural model of binocular saccade planning and vergence control. *Adaptive Behavior*, 23(5):265–282, 2015.

[129] Li, H, Ma, Y, and Principe, J. Cognitive architecture for video games. In *2020 International Joint Conference on Neural Networks (IJCNN)*, pages 1–9. IEEE, 2020.

[130] Rao, RPN, Olshausen, BA, and Lewicki, MS. *Probabilistic models of the brain: Perception and neural function*. MIT Press, 2002.

[131] Lee, TS and Mumford, D. Hierarchical Bayesian inference in the visual cortex. *JOSA A*, 20(7):1434–1448, 2003.

[132] Doya, K. *Bayesian brain: Probabilistic approaches to neural coding*. MIT press, 2007.

[133] Seth, AK. The cybernetic bayesian brain. In *Open mind*. Open MIND. Frankfurt am Main: MIND Group, 2014.

[134] Friston, K. The history of the future of the Bayesian brain. *NeuroImage*, 62(2):1230–1233, 2012.

[135] Clark, A. Whatever next? Predictive brains, situated agents, and the future of cognitive science. *Behavioral and Brain Sciences*, 36(3):181–204, 2013.

[136] Gabhart, KM, Xiong, YS, and Bastos, AM. Predictive coding: A more cognitive process than we thought? *Trends in Cognitive Sciences*, 2025.

[137] Gordon, N, Koenig-Robert, R, Tsuchiya, N, Van Boxtel, JJ, and Hohwy, J. Neural markers of predictive coding under perceptual uncertainty revealed with hierarchical frequency tagging. *eLife*, 6:e22749, 2017.

[138] Walsh, KS, McGovern, DP, Clark, A, and O'Connell, RG. Evaluating the neurophysiological evidence for predictive processing as a model of perception. *Annals of the New York Academy of Sciences*, 1464(1):242, 2020.

[139] Caucheteux, C, Gramfort, A, and King, JR. Evidence of a predictive coding hierarchy in the human brain listening to speech. *Nature Human Behaviour*, pages 1–12, 2023.

[140] Garrido, MI, Kilner, JM, Stephan, KE, and Friston, KJ. The mismatch negativity: A review of underlying mechanisms. *Clinical Neurophysiology*, 120(3):453–463, 2009.

[141] Wacongne, C, Changeux, JP, and Dehaene, S. A neuronal model of predictive coding accounting for the mismatch negativity. *Journal of Neuroscience*, 32(11):3665–3678, 2012.

[142] Weilnhammer, V, Stuke, H, Hesselmann, G, Sterzer, P, and Schmack, K. A predictive coding account of bistable perception-a model-based fMRI study. *PLoS Computational Biology*, 13(5):e1005536, 2017.

[143] Watanabe, E, Kitaoka, A, Sakamoto, K, Yasugi, M, and Tanaka, K. Illusory motion reproduced by deep neural networks trained for prediction. *Frontiers in Psychology*, 9:340023, 2018.

[144] Spratling, MW. Predictive coding as a model of response properties in cortical area V1. *Journal of Neuroscience*, 30(9):3531–3543, 2010.

[145] Kogo, N and Trengove, C. Is predictive coding theory articulated enough to be testable?, 2015.

[146] Spratling, MW. Fitting predictive coding to the neurophysiological data. *Brain Research*, 1720:146313, 2019.

[147] Ficco, L, Mancuso, L, Manuello, J, Teneggi, A, Liloia, D, Duca, S, Costa, T, Kovacs, GZ, and Cauda, F. Disentangling predictive processing in the brain: A meta-analytic study in favour of a predictive network. *Scientific Reports*, 11(1):16258, 2021.

[148] Heilbron, M and Chait, M. Great expectations: Is there evidence for predictive coding in auditory cortex? *Neuroscience*, 389:54–73, 2018.

[149] Bastos, AM, Usrey, WM, Adams, RA, Mangun, GR, Fries, P, and Friston, KJ. Canonical microcircuits for predictive coding. *Neuron*, 76(4):695–711, 2012.

[150] Shipp, S. Neural elements for predictive coding. *Frontiers in Psychology*, 7:1792, 2016.

[151] Schiess, M, Urbanczik, R, and Senn, W. Somato-dendritic synaptic plasticity and error-backpropagation in active dendrites. *PLoS Computational Biology*, 12(2):e1004638, 2016.

[152] Mikulasch, FA, Rudelt, L, Wibral, M, and Priesemann, V. Where is the error? Hierarchical predictive coding through dendritic error computation. *Trends in Neurosciences*, 46(1):45–59, 2023.

[153] Bastos, AM, Lundqvist, M, Waite, AS, Kopell, N, and Miller, EK. Layer and rhythm specificity for predictive routing. *Proceedings of the National Academy of Sciences*, 117(49):31459–31469, 2020.

[154] Hohwy, J, Roepstorff, A, and Friston, K. Predictive coding explains binocular rivalry: An epistemological review. *Cognition*, 108(3), 2008.

[155] Auksztulewicz, R and Friston, K. Repetition suppression and its contextual determinants in predictive coding. *Cortex*, 80, 2016.

[156] Feldman, H and Friston, KJ. Attention, uncertainty, and free-energy. *Frontiers in Human Neuroscience*, 4:215, 2010.

[157] Kanai, R, Komura, Y, Shipp, S, and Friston, K. Cerebral hierarchies: Predictive processing, precision and the pulvinar. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 370(1668):20140169, 2015.

[158] Bogacz, R. A tutorial on the free-energy framework for modelling perception and learning. *Journal of Mathematical Psychology*, 76:198–211, 2017.

[159] Brown, H, Adams, RA, Parees, I, Edwards, M, and Friston, K. Active inference, sensory attenuation and illusions. *Cognitive processing*, 14:411–427, 2013.

[160] Moran, RJ, Campo, P, Symmonds, M, Stephan, KE, Dolan, RJ, and Friston, KJ. Free energy, precision and learning: the role of cholinergic neuromodulation. *Journal of Neuroscience*, 33(19):8227–8236, 2013.

[161] Parr, T and Friston, KJ. Attention or salience? *Current Opinion in Psychology*, 29:1–5, 2019.

[162] Barron, HC, Auksztulewicz, R, and Friston, K. Prediction and memory: A predictive coding account. *Progress in Neurobiology*, 192:101821, 2020.

[163] Olshausen, BA and Field, DJ. Sparse coding with an overcomplete basis set: A strategy employed by V1? *Vision Research*, 37(23):3311–3325, 1997.

[164] Kinghorn, PF, Millidge, B, and Buckley, CL. Preventing deterioration of classification accuracy in predictive coding networks. In *Active Inference: Third International Workshop, IWAI 2022, Grenoble, France, September 19, 2022, Revised Selected Papers*, pages 1–15. Springer, 2023.

[165] Ororbia, A. Spiking neural predictive coding for continually learning from data streams. *Neurocomputing*, 544:126292, 2023.

[166] Alonso, N and Neftci, E. Tightening the biological constraints on gradient-based predictive coding. In *International Conference on Neuromorphic Systems 2021*, pages 1–9, 2021.

[167] Ororbia, A. Contrastive-signal-dependent plasticity: Forward-forward learning of spiking neural systems. *arXiv preprint arXiv:2303.18187*, 2023.

[168] Cornford, J, Pogodin, R, Ghosh, A, Sheng, K, Bicknell, BA, Codol, O, Clark, BA, Lajoie, G, and Richards, BA. Brain-like learning with exponentiated gradients. *bioRxiv*, pages 2024–10, 2024.

[169] PyTorch. https://github.com/pytorch/pytorch.

[170] Tensorflow. https://github.com/tensorflow/tensorflow.

[171] ngc-learn. https://github.com/ago109/ngc-learn.

[172] pypc. https://github.com/infer-actively/pypc.

[173] predify. https://github.com/miladmozafari/predify.

[174] Legrand, N, Weber, L, Waade, PT, Daugaard, AHM, Khodadadi, M, Mikuš, N, and Mathys, C. pyhgf: A neural network library for predictive coding. *arXiv preprint arXiv:2410.09206*, 2024.

[175] LeCun, Y. Deep learning hardware: Past, present, and future. In *2019 IEEE International Solid- State Circuits Conference - (ISSCC)*, pages 12–19, 2019.

[176] Chua, L. Memristor-the missing circuit element. *IEEE Transactions on Circuit Theory*, 18(5):507–519, 1971.

[177] Strukov, DB, Snider, GS, Stewart, DR, and Williams, RS. The missing memristor found. *Nature*, 453(7191):80–83, 2008.

[178] Grollier, J, Querlioz, D, Camsari, K, Everschor-Sitte, K, Fukami, S, and Stiles, MD. Neuromorphic spintronics. *Nature Electronics*, 3(7):360–370, 2020.

[179] Zucchet, N and Sacramento, J. Beyond backpropagation: Bilevel optimization through implicit differentiation and equilibrium propagation. *Neural Computation*, 34(12):2309–2346, 2022.

[180] Kendall, J, Pantone, R, Manickavasagam, K, Bengio, Y, and Scellier, B. Training end-to-end analog neural networks with equilibrium propagation. *arXiv preprint arXiv:2006.01981*, 2020.

[181] Kumar, S, Wang, X, Strachan, JP, Yang, Y, and Lu, WD. Dynamical memristors for higher-complexity neuromorphic computing. *Nature Reviews Materials*, 7(7):575–591, 2022.

[182] Haselager, DR, Boers, SN, Jongsma, KR, Vinkers, CH, Broekman, ML, and Bredenoord, AL. Breeding brains? patients' and laymen's perspectives on cerebral organoids. *Regenerative Medicine*, 15(12):2351–2360, 2020.

[183] Smirnova, L and Hartung, T. Neuronal cultures playing pong: First steps toward advanced screening and biological computing. *Neuron*, 110(23):3855–3856, 2022.

[184] Smirnova, L, Caffo, BS, Gracias, DH, Huang, Q, Morales Pantoja, IE, Tang, B, Zack, DJ, Berlinicke, CA, Boyd, JL, Harris, TD, et al. Organoid intelligence (OI): The new frontier in biocomputing and intelligence-in-a-dish. *Frontiers in Science*, page 0, 2023.

[185] Shipp, S, Adams, RA, and Friston, KJ. Reflections on agranular architecture: Predictive coding in the motor cortex. *Trends in Neurosciences*, 36(12):706–716, 2013.

[186] Friston, K. The sentient organoid? *Frontiers in Science*, page 0, 2023.

[187] Ballard, DH, Rao, RPN, and Zhang, Z. A single-spike model of predictive coding. *Neurocomputing*, 32:17–23, 2000.

[188] Wasserman, L. Bayesian model selection and model averaging. *Journal of Mathematical Psychology*, 44(1):92–107, 2000.

[189] De La Rocha, J, Doiron, B, Shea-Brown, E, Josić, K, and Reyes, A. Correlation between neural spike trains increases with firing rate. *Nature*, 448(7155):802, 2007.

[190] Maass, W. Networks of spiking neurons: The third generation of neural network models. *Neural Networks*, 10(9):1659–1671, 1997.

[191] Furber, SB, Galluppi, F, Temple, S, and Plana, LA. The SpiNNaker project. *Proceedings of the IEEE*, 102(5):652–665, 2014.

[192] Merolla, PA, Arthur, JV, Alvarez-Icaza, R, Cassidy, AS, Sawada, J, Akopyan, F, Jackson, BL, Imam, N, Guo, C, Nakamura, Y, et al. A million spiking-neuron integrated circuit with a scalable communication network and interface. *Science*, 345(6197):668–673, 2014.

[193] Davies, M, Srinivasa, N, Lin, TH, Chinya, G, Cao, Y, Choday, SH, Dimou, G, Joshi, P, Imam, N, Jain, S, et al. Loihi: A neuromorphic manycore processor with on-chip learning. *IEEE Micro*, 38(1):82–99, 2018.

[194] Ballard, D, Rao, R, and Zhang, Z. A model of predictive coding based on spike timing. *University of Rochester Computer Science Robotics and Vision Technical Reports*, 1999.

[195] Bi, Gq and Poo, Mm. Synaptic modifications in cultured hippocampal neurons: Dependence on spike timing, synaptic strength, and postsynaptic cell type. *Journal of Neuroscience*, 18(24):10464–10472, 1998.

[196] Hodgkin, AL and Huxley, AF. A quantitative description of membrane current and its application to conduction and excitation in nerve. *The Journal of Physiology*, 117(4):500–544, 1952.

[197] Lan, M, Xiong, X, Jiang, Z, and Lou, Y. PC-SNN: Supervised learning with local Hebbian synaptic plasticity based on predictive coding in spiking neural networks. *arXiv preprint arXiv:2211.15386*, 2022.

[198] Mikulasch, FA, Rudelt, L, Wibral, M, and Priesemann, V. Dendritic predictive coding: A theory of cortical computation with spiking neurons. *arXiv preprint arXiv:2205.05303*, 2022.

[199] N'dri, AW, Gebhardt, W, Teulière, C, Zeldenrust, F, Rao, RP, Triesch, J, and Ororbia, A. Predictive coding with spiking neural networks: A survey. *arXiv preprint arXiv:2409.05386*, 2024.

[200] Frieder, S and Lukasiewicz, T. (Non-)Convergence results for predictive coding networks. In *Proceedings of the 39th International Conference on Machine Learning, ICML 2022, Baltimore, Maryland, USA, 17-23 July 2022*, Proceedings of Machine Learning Research, pages 6793–6810. PMLR, 2022.

[201] Vilalta, R and Drissi, Y. A perspective view and survey of meta-learning. *Artificial Intelligence Review*, 18:77–95, 2002.

[202] Santoro, A, Bartunov, S, Botvinick, M, Wierstra, D, and Lillicrap, T. Meta-learning with memory-augmented neural networks. In *International Conference on Machine Learning*, pages 1842–1850. PMLR, 2016.

[203] Hospedales, T, Antoniou, A, Micaelli, P, and Storkey, A. Meta-learning in neural networks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(9):5149–5169, 2021.

[204] Kavukcuoglu, K, Ranzato, M, and LeCun, Y. Fast inference in sparse coding algorithms with applications to object recognition. *arXiv preprint arXiv:1010.3467*, 2010.

[205] Anil Meera, A and Wisse, M. Dynamic expectation maximization algorithm for estimation of linear systems with colored noise. *Entropy*, 23(10):1306, 2021.

[206] Friston, K. Predictive coding, precision and synchrony. *Cognitive Neuroscience*, 3(3-4):238–239, 2012.

[207] Khan, ME and Nielsen, D. Fast yet simple natural-gradient descent for variational inference in complex models. In *2018 International Symposium on Information Theory and Its Applications (ISITA)*, pages 31–35. IEEE, 2018.

[208] Osawa, K, Tsuji, Y, Ueno, Y, Naruse, A, Foo, CS, and Yokota, R. Scalable and practical natural gradient for large-scale deep learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(1):404–415, 2020.

[209] Tang, D and Ranganath, R. The variational predictive natural gradient. In *International Conference on Machine Learning*, pages 6145–6154. PMLR, 2019.

[210] Ofner, A and Stober, S. Predprop: Bidirectional stochastic optimization with precision weighted predictive coding. *arXiv preprint arXiv:2111.08792*, 2021.

[211] Alonso, N, Krichmar, J, and Neftci, E. Understanding and improving optimization in predictive coding networks. *arXiv preprint arXiv:2305.13562*, 2023.

[212] Zahid, U, Guo, Q, Friston, K, and Fountas, Z. Curvature-sensitive predictive coding with approximate Laplace Monte Carlo. *arXiv preprint arXiv:2303.04976*, 2023.

[213] Srivastava, N, Hinton, G, Krizhevsky, A, Sutskever, I, and Salakhutdinov, R. Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.

[214] Ioffe, S and Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, pages 448–456. pmlr, 2015.

[215] Kingma, DP and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[216] Tieleman, T, Hinton, G, et al. Lecture 6.5-RMSprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural Networks for Machine Learning*, 4(2):26–31, 2012.

[217] Fukushima, K. Visual feature extraction by a multilayered network of analog threshold elements. *IEEE Transactions on Systems Science and Cybernetics*, 5(4):322–333, 1969.

[218] Pascanu, R, Mikolov, T, and Bengio, Y. On the difficulty of training recurrent neural networks. In *International Conference on Machine Learning*, pages 1310–1318, 2013.

[219] Song, Y and Ermon, S. Generative modeling by estimating gradients of the data distribution. *Advances in Neural Information Processing Systems*, 32, 2019.

[220] Tschantz, A, Koudahl, M, Linander, H, Da Costa, L, Heins, C, Beck, J, and Buckley, C. Bayesian predictive coding. *arXiv preprint arXiv:2503.24016*, 2025.

[221] Dauwels, J. On variational message passing on factor graphs. In *2007 IEEE International Symposium on Information Theory*, pages 2546–2550. IEEE, 2007.

[222] Tenenbaum, JB, Kemp, C, Griffiths, TL, and Goodman, ND. How to grow a mind: Statistics, structure, and abstraction. *Science*, 331(6022):1279–1285, 2011.

[223] Tervo, DGR, Tenenbaum, JB, and Gershman, SJ. Toward the neural implementation of structure learning. *Current Opinion in Neurobiology*, 37:99–105, 2016.

[224] Smith, R, Schwartenbeck, P, Parr, T, and Friston, KJ. An active inference approach to modeling structure learning: Concept learning as an example case. *Frontiers in Computational Neuroscience*, 14:41, 2020.

[225] Rutar, D, de Wolff, E, van Rooij, I, and Kwisthout, J. Structure learning in predictive processing needs revision. *Computational Brain & Behavior*, 5(2):234–243, 2022.

[226] Kuchling, F, Friston, K, Georgiev, G, and Levin, M. Morphogenesis as Bayesian inference: A variational approach to pattern formation and control in complex biological systems. *Physics of Life Reviews*, 33:88–108, 2020.

[227] Kiebel, SJ and Friston, KJ. Free energy and dendritic self-organization. *Frontiers in Systems Neuroscience*, 5:80, 2011.

[228] Friston, K, Stephan, K, Li, B, Daunizeau, J, et al. Generalised filtering. *Mathematical Problems in Engineering*, 2010, 2010.

[229] Stephan, KE, Penny, WD, Moran, RJ, den Ouden, HE, Daunizeau, J, and Friston, KJ. Ten simple rules for dynamic causal modeling. *Neuroimage*, 49(4):3099–3109, 2010.

[230] Hoeting, JA, Madigan, D, Raftery, AE, and Volinsky, CT. Bayesian model averaging: A tutorial (with comments by M. Clyde, David Draper and EI George, and a rejoinder by the authors. *Statistical Science*, 14(4):382–417, 1999.

[231] Gershman, SJ and Blei, DM. A tutorial on Bayesian nonparametric models. *Journal of Mathematical Psychology*, 56(1):1–12, 2012.

[232] Penny, WD. Comparing dynamic causal models using AIC, BIC and free energy. *Neuroimage*, 59(1):319–330, 2012.

[233] Friston, K and Penny, W. Post hoc Bayesian model selection. *Neuroimage*, 56(4):2089–2099, 2011.

[234] Teh, YW, Jordan, MI, Beal, MJ, and Blei, DM. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006.

[235] Goldwater, SJ. *Nonparametric Bayesian Models of Lexican Acquisition*. Brown University, 2007.

[236] Laird, JE, Lebiere, C, and Rosenbloom, PS. A standard model of the mind: Toward a common computational framework across artificial intelligence, cognitive science, neuroscience, and robotics. *AI Magazine*, 38(4):13–26, 2017.

[237] Schwartenbeck, P and Friston, K. Computational phenotyping in psychiatry: A worked example. *eNeuro*, 3(4), 2016.

[238] Kanerva, P. *Sparse distributed memory*. MIT Press, 1988.

[239] Levy, SD and Gayler, R. Vector symbolic architectures: A new building material for artificial general intelligence. In *Proceedings of the 2008 Conference on Artificial General Intelligence 2008: Proceedings of the 1st AGI Conference*, pages 414–418, 2008.

[240] Kanerva, P. Hyperdimensional computing: An introduction to computing in distributed representation with high-dimensional random vectors. *Cognitive Computation*, 1:139–159, 2009.

[241] Ma, WJ and Peters, B. A neural network walks into a lab: Towards using deep nets as models for human behavior. *arXiv preprint arXiv:2005.02181*, 2020.

[242] Friston, K, FitzGerald, T, Rigoli, F, Schwartenbeck, P, Pezzulo, G, et al. Active inference and learning. *Neuroscience & Biobehavioral Reviews*, 68:862–879, 2016.

[243] Lotter, W, Kreiman, G, and Cox, D. Deep predictive coding networks for video prediction and unsupervised learning. *arXiv:1605.08104*, 2016.

[244] Sato, R, Kashima, H, and Yamamoto, T. Short-term precipitation prediction with skip-connected prednet. In *International Conference on Artificial Neural Networks*. Springer, 2018.

[245] Zhong, J, Cangelosi, A, Zhang, X, and Ogata, T. AFA-PredNet: the action modulation within predictive coding. In *2018 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2018.

[246] Rane, RP, Szügyi, E, Saxena, V, Ofner, A, and Stober, S. Prednet and predictive coding: A critical review. In *Proceedings of the 2020 International Conference on Multimedia Retrieval*, 2020.